

Validation of spatial variability in downscaling results from the VALUE perfect predictor experiment

M. Widmann¹, J. Bedia^{2,3}, J.M. Gutiérrez⁴, T. Bosshard⁵, E. Hertig⁶, D. Maraun⁷, M.J. Casado⁸, P. Ramos⁸, R.M. Cardoso⁹, P.M.M. Soares⁹, J. Ribalaygua¹⁰, C. Pagé¹¹, A. Fischer¹², S. Herrera², and R. Huth¹³

¹School of Geography, Earth and Environmental Sciences,
University of Birmingham, UK

²Dept. Applied Mathematics and Computing Science, University
of Cantabria, Santander, Spain

³Predictia Intelligent Data Solutions S.L., Santander, Spain

⁴National Research Council (CSIC), Instituto de Física de
Cantabria, Santander, Spain

⁵Swedish Meteorological and Hydrological Institute (SMHI),
Norrköping, Sweden

⁶Dept. of Geography, University of Augsburg, Germany

⁷Wegener Center for Climate and Global Change, University Graz,
Austria

⁸Agencia Estatal de Meteorología (AEMET), Madrid, Spain

⁹Instituto Dom Luiz (IDL), Faculdade de Ciências, Universidade
de Lisboa, Portugal

¹⁰Fundación para la Investigación del Clima (FIC), Spain

¹¹Centre Européen de Recherche et de Formation Avancée en
Calcul Scientifique (CERFACS), Toulouse, France

¹²Federal Office of Meteorology and Climatology (MeteoSwiss),
Zurich, Switzerland

¹³Institute of Atmospheric Physics, Charles University, Prague,
Czech Republic

February 5, 2019

Contents

1	Introduction	3
2	Data and methods	7
2.1	Observations and downscaled data	7
2.2	Validation measures	12
2.2.1	Correlations	12
2.2.2	Spatial degrees of freedom	12
2.2.3	Joint threshold exceedances	14
2.2.4	Regionalisation	15
3	Results	17
3.1	Example situation	17
3.2	Correlations	18
3.3	Spatial degrees of freedom	26
3.4	Joint threshold exceedances	28
3.5	Regionalisation	31
4	Summary and conclusions	35

Abstract

The spatial dependence of meteorological variables is crucial for many impacts, e.g. droughts, floods, river flows, energy demand, and crop yield. There is thus a need to understand how well it is represented in downscaling products. Within the COST Action VALUE we have conducted a comprehensive analysis of spatial variability in the output of over 40 different downscaling methods in a perfect predictor setup. The downscaling output is evaluated against daily precipitation and temperature observations for the period 1979-2008 at 86 sites across Europe and 53 sites across Germany. We have analysed the dependency of correlations of daily temperature and precipitation series at station pairs on the distance between the stations. For the European dataset we have also investigated the complexity of the downscaled data by calculating the number of independent spatial degrees of freedom. For daily precipitation at the German network we have additionally evaluated the dependency of the joint exceedance of the wet day threshold and of the local 90th percentile on the distance between the stations. Finally we have investigated regional patterns of European monthly precipitation obtained from rotated principal component analysis.

We analysed Perfect Prog methods, which are based on statistical relationships derived from observations, as well as Model Output Statistics approaches, which attempt to correct simulated variables. In summary we found that most Perfect Prog downscaling methods, with the exception of multi-site analog methods and a method that explicitly models spatial dependence yield unrealistic spatial characteristics. RCM-based Model Output Statistics methods showed good performance with respect to correlation lengths and the joint occurrence of wet days, but a substantial overestimation of the joint occurrence of heavy precipitation events. These findings apply to the spatial scales that are resolved by our observation network, and similar studies with higher resolutions, which are relevant for small hydrological catchment, are desirable.

1 Introduction

2 Projections for future climate change are primarily based on simulations with
3 coupled atmosphere-ocean general circulation models (GCMs). Their relatively
4 coarse horizontal resolution of around 100 km means that not all relevant at-
5 mospheric processes can be realistically modelled, which leads to errors on the
6 resolved scales. Moreover, the output does not have the spatial resolution often
7 needed for impact and adaptation studies. In order to overcome these problems
8 downscaling (DS) methods are routinely used, either based on high-resolution
9 regional climate models (RCMs), on statistical methods, or on a combination
10 of both (Maraun and Widmann, 2018; Ekstroem et al., 2015; Hewitson et al.,
11 2014; Maraun et al., 2010).

12 The spatial structure of the output from DS methods is highly relevant when
13 the results are used to assess impacts that are determined by spatial aggregation
14 of meteorological variables. Typical examples for which a realistic representa-
15 tion of spatial variability matters are river flow and floods (Arnaud et al., 2002;
16 Segond et al., 2007; Viviroli et al., 2009), droughts (Trambauer et al., 2015),
17 glacier mass balance (Machguth et al., 2009), ecosystem composition (Mon-
18 estiez et al., 2001), crop yields (Holzkämper et al., 2012), energy consumption
19 and production, as well as weather-related health problems. For instance an
20 over- or underestimation of correlations between precipitation timeseries at dif-
21 ferent locations within a river catchment would typically lead to an over- or
22 underestimation of high and low river flow conditions.

23 Within the COST action VALUE a comprehensive validation framework for
24 DS methods has been designed and implemented (Maraun et al., 2015). The
25 user-relevant aspects of DS output identified in the framework are marginal
26 distributions including extremes, temporal variability, and intervariable rela-
27 tionships, all considered at individual locations, as well as spatial variability.
28 The performance of DS methods with respect to the aspects defined at indi-
29 vidual stations within Europe has been investigated in the companion papers
30 in this special issue (Gutiérrez et al., 2018; Hertig et al., 2018; Maraun et al.,
31 2018). Here we analyse specifically how well the different DS methods represent
32 the spatial structure of precipitation and temperature fields over Europe. As
33 pointed out in Maraun et al. (2015) it is usually not the spatial pattern of the
34 long-term mean but the structure of the individual events that is relevant for
35 impacts, because it includes for instance the information on whether all loca-
36 tions within a river catchment tend to receive precipitation at the same time,
37 or whether it is likely that some areas stay dry when there is precipitation in
38 others. It can be useful to remove the effect of the climatological mean on indi-
39 vidual events and to analyse the residual spatial variability, i.e. to express the
40 data as deviations from the long-term mean.

41 More formally speaking, when considering a meteorological variable simul-
42 taneously at different locations we are dealing with a multivariate dataset given
43 by the values at the different locations, and the goal when validating spatial
44 variability is to investigate the similarity of the observed and downscaled data
45 clouds. To a first order approximation the datasets are characterised by their
46 multivariate long-term temporal means, i.e. by the patterns of the climatologi-
47 cal mean. For the observations it is mainly influenced by the meridional gradient
48 and local differences in the radiation budget, the proximity to the oceans, the
49 mean large-scale atmospheric circulation, and topography. These factors in-

50 fluence meteorological processes such as atmospheric stability, convection, flow
51 convergence, frontal passages, or Foehn, which affect the spatial structure of
52 individual weather events as well as of the long-term mean. It can be expected
53 that almost all statistical DS method reproduce the mean temperature and pre-
54 cipitation fields quite well by construction, for instance by estimating anomalies
55 around the observed mean in the case of regression-based methods or by ad-
56 justing distributions. The skill of DS methods with respect to representing the
57 mean has been analysed to some extent in Gutiérrez et al. (2018), albeit without
58 explicitly investigating the spatial pattern of the bias of the long-term means.
59 The mean bias in the raw output of regional models has been investigated in
60 many publications (e.g. Kotlarski et al., 2014; Isotta et al., 2015). Moreover,
61 as already mentioned, it is mostly the structure of the residual spatial variabil-
62 ity that is impact-relevant. We therefore focus in our analysis on the spatial
63 structure of the residual variability, mainly on the daily timescale.

64 For multivariate Gaussian data the structure of the variability around the
65 mean is fully captured by the covariance matrix, and for normalised data by
66 the correlation matrix. It is thus a natural starting point to investigate the
67 similarity of the observed and the downscaled covariances or correlations be-
68 tween different locations. As correlations are a direct measure for the strength
69 of linear relationships between timeseries we will consider those. We will also
70 investigate the probabilities for joint exceedances of thresholds, which are of
71 practical relevance for impact modelling and which for non-Gaussian data do
72 not directly follow from the covariance matrix. We note that multivariate data
73 can alternatively be described by a combination of their marginal distributions,
74 which are investigated in Gutiérrez et al. (2018), and copulas that analytically
75 express the dependence structure. However, for brevity this approach is not
76 taken here. In addition we will analyse the overall complexity, and the repre-
77 sentation of regional patterns. Details on our validation approach are given in
78 the method section.

79 In spite of the importance of the spatial structure of daily values for cli-
80 mate impacts, only a few studies have validated the spatial aspects of stan-
81 dard deterministic Perfect Prog (PP) downscaling products. Correlations be-
82 tween timeseries at different locations, including their dependency on distance,
83 have been analysed (Easterling, 1999; Kettle and Thompson, 2004; Huth et al.,
84 2008, 2015), and homogeneous regions have been investigated by cluster analy-
85 sis (Huth, 2002). These studies, most of which focus on temperature, indicate
86 that PP methods that use large-scale predictors overestimate spatial correla-
87 tions, whereas local analog methods underestimate them. Huth et al. (2015)
88 additionally included two RCMs in the method comparison and found no sys-
89 tematic over- or underestimation for them. A comparison of some PP and MOS
90 methods, as well as RCMs, undertaken by Ayar et al. (2016) included some
91 analysis of spatial variability of daily precipitation based on the leading Prin-
92 cipal Component (PC) loading patterns and on correlations of daily patterns.
93 The study found a mixed performance of the RCMs and MOS with better skill
94 in winter than in summer, and in general low performance for PP methods. The
95 analog method showed as expected realistic PC loadings but failed to capture
96 the individual daily patterns.

97 In addition, stochastic PP methods that explicitly model spatial structure
98 have been developed and analysed. Frost et al. (2011) evaluated correlations
99 of occurrence and amount of daily precipitation at different locations obtained

100 from a Nonhomogeneous Hidden Markov Model (NHMM) for occurrence com-
101 bined with conditional multiple regression for amounts, and from GLIMCLIM, a
102 conditional multisite weather generator based on a generalised linear model, and
103 found that both substantially underestimated intersite correlations. Hu et al.
104 (2013) obtained similar results for GLIMCLIM, but found in contrast that a
105 NHMM performed well. The difference can be a result of both the predictor
106 choice, or the specific regional climate. A further method type are conditional
107 multisite weather generators for precipitation constrained by the observed de-
108 pendences between sites, which were found to represent the observed properties
109 well (Cannon, 2008; Wilks, 2012).

110 Disaggregation methods for precipitation investigated in Ferraris et al.
111 (2003) show substantial over- and underestimations of intersite correlations with
112 no method performing systematically better than others. However, advanced
113 stochastic models for precipitation that include a disaggregation step based on
114 two-dimensional, latent Gaussian fields showed realistic spatial characteristics
115 (Paschalis et al., 2013).

116 Recently several analog methods in which the analogs are based on a coarse
117 resolution representation of the predictand variable rather than on the large-
118 scale atmospheric circulation have been developed. There are different imple-
119 mentations depending on how model biases are treated and on how the down-
120 scaled field is constructed from a pool of analog situations; for a description of
121 the frequently used 'localised constructed analog method' (LOCA) and a dis-
122 cussion of other variants see Pierce et al. (2014). They are implemented such
123 that a common analog is chosen for adjacent locations and thus yield realistic
124 spatial fields by construction if individual analogs are used and fairly realistic
125 fields if weighted means of multiple analogs are used. An intercomparison of bias
126 corrected constructed analogs (BCCA), of methods combining bias correction
127 for monthly or daily fields and spatial disaggregation (BCSDm, BCSMd), and
128 of an asynchronous regression method is presented in Gutmann et al. (2014),
129 who found that all methods but BCSDm substantially overestimate spatial cor-
130 relations. The reason for the good performance of BCSDm is that in contrast
131 to the other methods it inherits the spatial variability from the observations,
132 rather than from the driving model.

133 Recent developments also include multisite MOS methods. Bárdossy and
134 Pegram (2012) found that RCM precipitation had too low intersite correlations
135 and formulated a matrix and a sequential recorrelation method to adjust the
136 spatial structure, with the former applicable to match Pearson correlations and
137 the latter to reproduce more general copula-based representations of the mul-
138 tivariate structure. The correction methods led to a realistic spatial structure,
139 with the exception of an underrepresentation of clustering of extreme precipita-
140 tion, allow for changes in the spatial dependences in a future climate, and mainly
141 preserve the temporal structure of the RCM output. Cannon (2018) developed
142 a multivariate quantile mapping method that yields the observed multivariate
143 distribution, applied it to correct spatial RCM precipitation fields, and demon-
144 strated realistic spatial characteristics of the corrected fields. There are also
145 parametric quantile mapping methods that interpolate the observed distribu-
146 tion parameters to high spatial resolution (Mamalakis et al., 2017), but as they
147 do not model the spatial structure of variability they are essentially singlesite
148 MOS methods.

149 In the context of ensemble weather forecasting postprocessing methods have

150 been used that rearrange the simulated data in time so they have the same rank
151 structure as the observations in a training period (known as Schaake Shuffle),
152 which leads to a reproduction of the spatial and intervariable dependence struc-
153 ture of the training data (Clark et al., 2004). The method has been employed to
154 provide input for hydrological forecasts (Voisin et al., 2011) and to postprocess
155 atmospheric reanalyses (Vrac and Friederichs, 2015). A drawback that makes its
156 application in a climate change context problematic is that it is constrained to
157 reproduce the temporal rank structure of the training dataset. Vrac (2018) has
158 suggested a rank-based resampling method that relaxes this condition and also
159 introduces stochasticity by generating as many multivariate corrected outputs
160 as the number of statistical dimensions (i.e. number of grid-cells \times number of
161 climate variables). This study has also demonstrated how to apply the method
162 in a climate change context. However, further research on the usefulness of the
163 method for climate change studies is needed, for instance because the reshuffling
164 breaks the physical consistency between large-scale atmospheric states and the
165 postprocessed variables, and will usually modify the climate change signal.

166 Our analysis extends these studies by considering a large number of down-
167 scaling methods (47 for precipitation and 45 for temperature) and by systemati-
168 cally comparing them with respect to several measures of spatial variability, us-
169 ing validation datasets over Europe and Germany. The structure of DS methods
170 can be expected to have a strong influence on the spatial variability of their out-
171 put. Singlesite methods, which are fitted to individual target locations, might
172 for instance yield a realistic spatial structure if the predictors explain a large
173 fraction of the local variability, but might overestimate spatial correlations if
174 small-scale variability is substantial and not adequately represented. A detailed
175 analysis of the variance explained by each downscaling method is provided in
176 Gutiérrez et al. (2018). Multisite DS methods, which simultaneously use sev-
177 eral locations for model fitting, might either achieve realistic spatial variability
178 through the common influence of predictors or through explicit constraints on
179 the multivariate structure of noise components or of the final output. In our
180 study we compare downscaling methods of different types which will allow us
181 to investigate whether some types exhibit a common behaviour with respect to
182 spatial variability. We note that the VALUE perfect predictor experiment uses
183 an ensemble of opportunity in which most of the methods are fitted on single
184 sites, reflecting the dominance of such methods in DS applications. In par-
185 ticular, no method explicitly models spatial dependence in the European-wide
186 experiment, although for some methods, spatial dependence results as a conse-
187 quence of the use of common predictors (e.g. regression methods using PCs) or
188 of the method characteristics (e.g. some analog methods using the same analog
189 day for all sites). However, for the additional experiment over Germany, two
190 regression methods that explicitly consider spatial dependence have contributed
191 to the study.

192 Section 2 starts with a discussion of the observations used for validation
193 as well as of the downscaled data, including a brief overview of the different
194 types of downscaling methods and of the experimental setup. It then continues
195 with an explanation of the different measures for spatial variability employed
196 to validate and compare the downscaling methods. Section 3 will present the
197 validation results in separate subsections for each validation measure. Summary
198 and conclusions will be given in section 4.

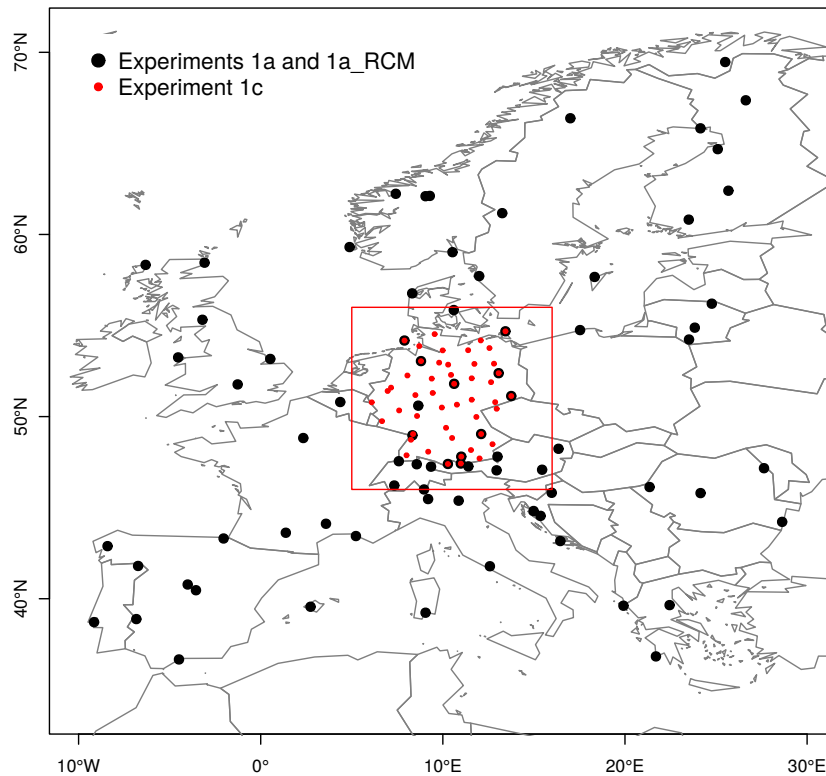


Figure 1: Locations of the reference stations for the European experiments (1a and 1a-RCM, black circles, VALUE-ECA-86-v2 dataset) and the German experiment (1c, red, VALUE-ECA-53-Germany-spatial-v1 dataset).

199 2 Data and methods

200 2.1 Observations and downscaled data

201 The predictands for the DS methods are observations for daily precipitation
 202 as well as for daily minimum and maximum temperature at 86 stations across
 203 Europe. This VALUE-ECA-86-v2 dataset is a subset of the publicly available
 204 ECA dataset (Tank et al., 2002) and covers the period 1979 - 2008. Besides the
 205 European-wide experiment (referred to as experiment_1a, or simply exp_1a),
 206 which is the common experiment for the different validation studies, we also
 207 present here the results of an experiment based on a denser ECA subset of 53
 208 stations within Germany for the same variables (referred to as experiment_1c,
 209 or simply exp_1c), which was designed to focus on spatial validation aspects.
 210 Details on data availability are given in Gutiérrez et al. (2018). Both networks
 211 are shown in Fig. 1.

212 The downscaling methods that have been considered in our study for precip-
 213 itation are listed in Table 1, those used for temperature in Table 2. The columns
 214 ‘1a’ and ‘1c’ indicate the methods contributing to each of the experiments. All
 215 downscaling methods have been calibrated following a five-fold cross validation
 216 with non-overlapping consecutive 6-year blocks. Further details about the meth-
 217 ods and the experimental setup can be found in Maraun et al. (2015), Gutiérrez

218 et al. (2018), and on www.value-cost.eu/validation#Experiment_1a.

219 We distinguish between PP and MOS methods (see e.g. Maraun et al.,
220 2010)). For the former the statistical relationships are derived from observa-
221 tions whereas MOS methods are fitted using predictors from RCMs (or global
222 climate models). PP methods represent real-world links between large-scale
223 predictors and the local predictand, and thus in applications to output from
224 climate models they require realistically simulated predictors – hence the name
225 'Perfect Prog(nosis)'. MOS methods represent relationships between simulated
226 and observed variables, are therefore model-specific, and do not only repre-
227 sent downscaling relationships but can also correct model biases. Unconditional
228 weather generators (WGs), which are statistical models that produce timeseries
229 with temporal characteristics similar to observations without any predictors are
230 a third group of methods listed under 'WG'. Conditional WGs, which include
231 meteorological predictors that influence the properties of the timeseries, should
232 not be categorised as a separate group to MOS and PP, because depending on
233 the setup for model fitting they either follow the PP or MOS approach, and are
234 thus listed under either PP or MOS.

235 The PP methods are validated in a perfect predictor setup using predictors
236 from the ERA-Interim Reanalysis (Dee et al., 2011) for the period 1979 - 2008 on
237 a coarse-grained 2° resolution, which is similar to typical output from global cli-
238 mate models. The PP assumption for the predictors is thus met by construction.
239 The MOS methods for the European experiment exp_1a are directly applied to
240 the ERA-Interim data on both the original 0.75° and on the coarse-grained res-
241 olution. We have conducted an additional European experiment exp_1a_RCM
242 for which the MOS predictors are taken from the RACMO RCM (van Meijgaard
243 et al., 2008) driven by perfect boundary conditions from ERA-Interim on the
244 original 0.75° resolution. For the German experiment exp_1c we have used MOS
245 predictors from ERA-Interim on the original 0.75° resolution.

246 The PP methods used here cover the widely used approaches, i.e. analog,
247 regression and weather type methods; the MOS methods cover frequently used
248 quantile mapping methods as well as recently developed stochastic MOS.

249 Information on the structural elements of the DS methods that may influence
250 the spatial characteristics of the output are also given in tables 1 and 2. The
251 'MS' column indicates whether the DS model has been fitted simultaneously for
252 multiple (or all) locations ('yes') or individually for each location ('no'). The
253 'EX' column lists whether the statistical model has explicit constraints on the
254 structure of spatial variability ('yes'), for instance on correlations for adjacent
255 locations. The 'ST' column indicates whether the DS output contains stochastic
256 noise ('yes'). The final column 'PC' states whether or not principal components
257 have been used as predictors. As already mentioned almost all of the methods
258 are fitted and applied at single sites, with only some analog methods being
259 applied to multiple sites. Note that methods that are fitted at individual sites
260 might still be used for multiple sites if for instance realistic spatial patterns can
261 be expected through the influence of the predictors.

262 All methods participating in the European experiment are fully described in
263 Annex 1 of Gutiérrez et al. (2018). We now describe the two additional methods,
264 GLM-BN-DET and DSCLIM-D, contributing only to the German experiment.
265 GLM-BN-DET is a multivariate extension of the GLM-DET method, which
266 explicitly models the spatial structure of precipitation occurrence by considering
267 a dependence graph linking marginally and/or conditionally dependent stations.

268 This graph allows to obtain a probabilistic model (a Bayesian network) which
269 encodes all the dependences displayed in the graph by means of an appropriated
270 factorisation of the joint probability distribution. This model allows simulating
271 spatially consistent precipitation occurrences. Moreover, for each particular
272 station, the model determines the set of stations (Markov blanket) exerting
273 a spatial influence. For each station, this set is included as spatial predictors
274 (in addition to the large-scale information) in the binomial/gamma GLM model
275 thus the model yields spatially consistent precipitation amounts. Details on this
276 particular methodology are given in Cano et al. (2004). DSCLIM-D is based
277 on weather typing, combined with linear regression and weather analogs. The
278 method has been introduced by Boé et al. (2006), but the version used here
279 differs in some details. The implementations for temperature and precipitation
280 are slightly different, and for brevity we explain only the latter case. DSCLIM-D
281 uses a clustering method to determine weather types (10 in this implementation)
282 in the SLP field. For each day the Euclidean distances of the SLP field to all the
283 weather types are calculated and used as predictors for the square root of the
284 precipitation anomaly at a given location in a multiple linear regression. The
285 mean of the estimated precipitation over all stations in the target area is then
286 used to define a set of analog days from which the downscaled local precipitation
287 is chosen. The set is defined by the days in the fitting period that belong to the
288 same weather type as well as have averaged precipitation in the same decile as
289 the estimated averaged precipitation. We note that comparing deciles is similar
290 to quantile mapping or inflated regression. In the deterministic version of the
291 method, which is used here, one analog precipitation field is randomly selected,
292 the stochastic version used several analogs.

Type	Code	Tech	1a	1c	MS	EX	ST	PC
MOS	Ratyetal-M6	S	×	-	no	no	no	no
	Ratyetal-M7	S	×	-	no	no	no	no
	ISI-MIP	S/PM	×	×	no	no	no	no
	DBS	PM	×	×	no	no	no	no
	Ratyetal-M9	PM	×	-	no	no	no	no
	BC	PM	×	×	no	no	no	no
	GQM	PM	×	×	no	no	no	no
	GPQM	PM	×	×	no	no	no	no
	EQM	QM	×	×	no	no	no	no
	EQMs	QM	×	-	no	no	no	no
	EQM-WT	QM/WT	×	×	no	no	no	no
	QMm	QM	×	×	no	no	no	no
	QMBC-BJ-PR	QM	×	-	no	no	no	no
	CDFt	QM	×	-	no	no	no	no
	QM-DAP	QM	×	-	no	no	no	no
	EQM-WIC658	QM	×	-	no	no	no	no
	Ratyetal-M8	QM	×	-	no	no	no	no
	MOS-AN	A	×	-	yes	no	no	no
	MOS-GLM	TF	×	-	no	no	yes	no
	VGLMGAMMA	TF/WG	×	-	no	no	yes	no
FIC02P	PM/A/TF	×	×	no	no	no	no	
FIC04P	PM/A/TF	×	×	no	no	no	no	
PP	FIC01P	A/TF	×	×	yes	no	no	no
	FIC03P	A/TF	×	×	yes	no	no	no
	ANALOG-ANOM	A	×	-	yes	no	no	no
	ANALOG	A	×	×	yes	no	no	yes
	ANALOG-MP	A	×	×	yes	no	yes	no
	ANALOG-SP	A	×	-	yes	no	yes	no
	MO-GP	TF	×	-	no	no	no	no
	GLM-P	TF	×	×	no	no	yes ^(a)	no
	MLR-RAN	TF	×	×	no	no	no	no
	MLR-RSN	TF	×	×	no	no	no	no
	MLR-ASW	TF	×	-	no	no	yes	no
	MLR-ASI	TF	×	×	no	no	no	no
	GLM-DET	TF	×	×	no	no	no	yes
	GLM	TF	×	-	no	no	yes	yes
	GLM-WT	TF/WT	×	×	no	no	yes	yes
	GLM-BN-DET	TF	-	×	yes	yes	no	yes
	DSCLIM-D	A/WT	-	×	yes	no	no	no
WT-WG	WT	×	-	no	no	yes	yes	
SWG	TF	×	-	no	no	yes	yes	
WG	SS-WG	WG	×	-	no	no	yes	no
	MARFI-BASIC	WG	×	-	no	no	yes	no
	MARFI-TAD	WG	×	-	no	no	yes	no
	MARFI-M3	WG	×	-	no	no	yes	no
	GOMEZ-BASIC	WG	×	-	no	no	yes	no
	GOMEZ-TAD	WG	×	-	no	no	yes	no

Table 1: Participating methods for precipitation for the European (exp1a) and German experiment (exp1c). Techniques: A: analog; S: scaling; PM: parametric quantile mapping; QM: empirical quantile mapping; TF: regression-like transfer function; WT: weather typing; WG: weather generator. Columns 1a and 1c indicate whether the methods have participated in the European and German experiment. MS: Multisite fitting; MS; EX: Explicitly modelled spatial structure; ST: Stochastic noise; PC: PCs used as predictors. ^(a) Only occurrence is randomised, amounts are based on inflated regression (in this case, the results are based on a single realisation).

Type	Tech	Code	MS	EX	ST	PC
MOS	RaiRat-M6	S	no	no	no	no
	RaiRat-M7	S	no	no	no	no
	RaiRat-M8	S	no	no	no	no
	SB	S	no	no	no	no
	ISI-MIP	S/PM	no	no	no	no
	DBS	PM	no	no	no	no
	GPQM	PM	no	no	no	no
	EQM	QM	no	no	no	no
	EQMs	QM	no	no	no	no
	EQM-WT	QM/WT	no	no	no	no
	QMm	QM	no	no	no	no
	QMBC-BJ-PR	QM	no	no	no	no
	CDFt	QM	no	no	no	no
	QM-DAP	QM	no	no	no	no
	EQM-WIC658	QM	no	no	no	no
	RaiRat-M9	QM	no	no	no	no
	DBBC	QM	no	no	no	no
	DBD	QM	no	no	no	no
	MOS-REG	TF	no	no	no	no
	FIC02T	PM/A/TF	no	no	no	no
PP	FIC01T	A/TF	yes	no	no	no
	ANALOG-ANOM	A	yes	no	no	no
	ANALOG	A	yes	no	no	yes
	ANALOG-MP	A	yes	no	yes	no
	ANALOG-SP	A	yes	no	yes	no
	MO-GP	TF	no	no	no	no
	MLR-T	TF	no	no	no	no
	MLR-RAN	TF	no	no	no	no
	MLR-RSN	TF	no	no	no	no
	MLR-ASW	TF	no	no	yes	no
	MLR-ASI	TF	no	no	no	no
	MLR-AAN	TF	no	no	no	no
	MLR-AAI	TF	no	no	no	no
	MLR-AAW	TF	no	no	yes	no
	MLR-PCA-ZTR	TF	no	no	no	yes
	MLR	TF	no	no	no	yes
	MLR-WT	TF/WT	no	no	no	yes
	WT-WG	WT	no	no	yes	yes
SWG	TF	no	no	yes	yes	
WG	SS-WG	WG	no	no	yes	no
	MARFI-BASIC	WG	no	no	yes	no
	MARFI-TAD	WG	no	no	yes	no
	MARFI-M3	WG	no	no	yes	no
	GOMEZ-BASIC	WG	no	no	yes	no
	GOMEZ-TAD	WG	no	no	yes	no

Table 2: Participating methods for temperature for the European experiment (exp1a). Techniques: A: analog; S: scaling; PM: parametric quantile mapping; QM: empirical quantile mapping; TF: regression-like transfer function; WT: weather typing; WG: weather generator. Multisite fitting: MS; EX: Explicitly modelled spatial structure; ST: Stochastic noise; PC: PCs used as predictors.

293 2.2 Validation measures

294 We now discuss the different validation measures on which the method compar-
295 ison is based. All computations have been done in R and the codes are publicly
296 available at Santander Meteorology Group (2016).

297 2.2.1 Correlations

298 Pairwise cross-correlations among all pairs of stations ($n \times \frac{n-1}{2}$ pairs, n being
299 the number of stations) are computed for the different target variables and
300 seasons (Spearman for precipitation and Pearson for temperatures), and for
301 experiments 1a and 1a-RCM ($n = 86$) and 1c ($n = 53$). For the temperature
302 data the seasonal cycle of each data series is removed prior to correlation analysis
303 by subtracting the climatological mean for each particular day of the year based
304 on the whole analysis period 1979-2008. The mean is based on a circular moving
305 average with a window width of 31 days centred around the target day. The
306 precipitation data are used in their original form. In both cases, no detrending
307 has been used. In addition to the visual comparison of correlation matrices
308 we calculate the correlation matrix distance (CMD, Herdin et al., 2005). It
309 measures the similarity between two correlation matrices and is defined as one
310 minus the inner product of the normalised vectorized matrices. For matrices
311 that are identical up to a scaling factor, the CMD is zero and for very different
312 matrices, for which the associated vectors are orthogonal, the CMD is one.

313 Station correlograms are then derived by plotting the cross-correlation value
314 for each station pair against their respective (great circle) geographical dis-
315 tances. As the resulting cloud of points may hinder a quick assessment of the
316 dependency of the correlations on distance, we fitted reference curves to each
317 correlogram using a local polynomial fit (“loess”, degree 2), allowing for a bet-
318 ter comparability between downscaling methods and against the reference data.
319 The local fit was preferred to other correlogram global fitting models commonly
320 used in geostatistics (e.g. exponential or spherical; see e.g. Hengl, 2007)), as it
321 does not require *a priori* assumptions about the structure of the correlations.
322 It is therefore suitable for different kinds of correlation structures and flexible
323 enough to allow for a direct comparison across different downscaling methods
324 and experiments. As a measure for the overall behaviour of the fitted curves we
325 then calculated correlation lengths (CL) for certain representative thresholds,
326 as the abscissa of the point of intersection of the correlation threshold with the
327 fitted line. We tested different thresholds, and the final values used are given in
328 Table 3. The CL biases for the predictions were calculated as the difference of
329 the CL for a given method and the CL of the observations (Table 4). This bias
330 is a simple measure for the difference in the correlation structure between the
331 predictions and the observations.

332 2.2.2 Spatial degrees of freedom

333 We determine the number of independent spatial degrees of freedom (DOF) that
334 are associated with the observations and with the downscaling products. DOFs
335 quantify the complexity of time- and space-dependent datasets and are based on
336 the correlation or covariance matrix. In addition to describing the dependency

Var.	Exps. 1a and 1a-RCM	Exp. 1c
Precip	0.35	0.50
Tmin	0.50	0.65
Tmax	0.50	0.65

Table 3: Correlation thresholds used for calculating correlation lengths in the European experiments (1a and 1a-RCM) and in the German experiment (1c).

Var.	Exps. 1a and 1a-RCM					Exp. 1c				
	<i>annual</i>	<i>DJF</i>	<i>JJA</i>	<i>MAM</i>	<i>SON</i>	<i>annual</i>	<i>DJF</i>	<i>JJA</i>	<i>MAM</i>	<i>SON</i>
Precip	495	527	429	475	540	404	546	310	393	417
Tmin	741	822	647	771	653	569	695	462	541	475
Tmax	873	1005	785	893	870	698	788	697	653	668

Table 4: Correlation length (CL) values (in km) calculated from the correlograms of the reference station datasets (VALUE-ECA-86-v2 for experiments 1a and 1a-RCM, and VALUE-ECA-53-Germany-spatial-v1 for experiment 1c).

337 of the correlations on distance by a single number (CL) we thus also use a single
 338 number to capture a key property of the correlation matrices themselves and
 339 then calculate its biases.

340 One possible way to define complexity is to consider the eigenvalue spec-
 341 trum of the covariance or correlation matrix. Consider a situation where the
 342 timeseries at all locations are perfectly correlated, which means there would be
 343 only one independent variable. In this case one PC (e.g. Hannachi et al., 2007)
 344 would explain all the variance, i.e. the first eigenvalue of the covariance matrix
 345 would be equal to the total variance and all other eigenvalues would be zero.
 346 If, in the other extreme case, the timeseries at all locations were independent,
 347 the eigenvalue spectrum would be completely flat, as no correlations between
 348 the station records could be exploited to construct any PCs that explain more
 349 variance than an individual station record. Roughly speaking, the steepness of
 350 the eigenvalue spectrum can thus be taken as an indication for the complex-
 351 ity of the data, with a steep (flat) spectrum being associated with low (high)
 352 complexity.

353 An alternative way to define the complexity of a space- and time-dependent
 354 field $\psi_i(t)$ is to consider the timeseries of the spatial sum of the squares of the
 355 values at the individual locations i , i.e.

$$E(t) = \sum_{i=1}^n \psi_i^2(t) \quad (1)$$

356 with n being the number of locations. For independent variables $E(t)$ has a χ^2 -
 357 distribution with N degrees of freedom, for dependent variables the distribution
 358 is well approximated by a χ^2 distribution with fewer degrees of freedom. A useful
 359 measure of complexity is obtained by asking how many independent variables
 360 are needed to obtain approximately the same χ^2 distribution, which is defined
 361 by its mean and variance, as for the timeseries of the sum of squares of the
 362 dependent variables.

363 This approach has been reviewed by Bretherton et al. (1999) who have shown

364 that for normally-distributed PCs the χ^2 and the eigenvalue approaches are
 365 equivalent if, as suggested in earlier studies, the degrees of freedom (DOF) are
 366 calculated from the eigenvalue spectrum by:

$$DOF = \frac{(\sum_i \lambda_i)^2}{\sum_i \lambda_i^2} \quad (2)$$

367 where λ_i is the i -th eigenvalue and the summation is over all the eigenvalues.

368 In this paper we follow the computationally easier eigenvalue approach and
 369 calculate the independent spatial degrees of freedom according to equation 2.
 370 The normality assumption has been checked in the reference observation dataset
 371 VALUE-ECA-86-v2 (see Sec. 2.1), by comparing the empirical distribution func-
 372 tion of each PC against the cumulative distribution function of the normal
 373 distribution using the Kolmogorov-Smirnov test, implemented in the function
 374 `ks.test` of the R package `stats` (R Core Team, 2018). All PCs were found
 375 to be indistinguishable from a normal distribution at the 5% significance level.
 376 The singular value decomposition implementation used in the (function `svd`
 377 in the R package `stats` R Core Team, 2018)) cannot handle missing values in
 378 the covariance matrix, and a few methods yielding missing values for all data
 379 in some stations did thus not yield results (this will be later indicated in the
 380 corresponding figure captions).

381 For consistency with the analysis of correlation lengths (Sec. 2.2.1), we base
 382 the DOFs on the eigenvalues of the correlation rather than the covariance ma-
 383 trix. In other words, we calculate the DOFs for standardised data, where the
 384 timeseries at each location have the same variance. The seasonal cycle is sub-
 385 tracted in the same way as for the correlation analysis. The DOFs for the
 386 observations, which are the reference for calculating DOF biases, are given in
 387 Table 5.

	DJF	MAM	JJA	SON
precip	30.02	41.51	48.64	36.05
tmin	6.56	7.66	9.43	8.86
tmax	5.65	6.56	7.55	6.91

Table 5: Degrees of freedom (DOF) for daily precipitation, minimum and max-
 imum temperature from the VALUE-ECA-86-v2 observation dataset.

388 2.2.3 Joint threshold exceedances

389 The correlation-based analyses discussed above investigate the strength of lin-
 390 ear relationships between the timeseries at different locations. However, for
 391 users of the downscaled data it may often be also relevant to know whether the
 392 probabilities for joint exceedance of a certain threshold at different locations are
 393 realistic in the downscaled data. Typical examples are the joint occurrence of
 394 precipitation or of heavy precipitation. For brevity, we restrict the analysis of
 395 such joint threshold exceedances to precipitation. This is the most challenging
 396 case since temperature fields are typically much smoother and spatially homo-
 397 geneous. Therefore, we consider two typical cases: the wet day threshold of 1
 398 mm/day and exceedance of thresholds for high precipitation, namely the local
 399 90th percentile.

400 The most direct way to analyse the dependence between the data X_i, X_j at
 401 a pair of stations $\{i, j\}$ for exceeding a threshold x_{0i} at location i and x_{0j} at
 402 location j , is subtracting the product of marginals $P(x_i \geq x_{0i}) \cdot P(x_j \geq x_{0j})$
 403 from the joint probability $P(x_i \geq x_{0i}, x_j \geq x_{0j})$. Their difference is zero only in
 404 case that $P(x_i \geq x_{0i})$ and $P(x_j \geq x_{0j})$ are totally independent and the larger
 405 the value, the more dependent they are. However, this difference would not
 406 only be influenced by the dependence for threshold exceedance, but also by the
 407 marginal probabilities at each of the stations, and is thus not a useful measure
 408 for the dependence itself.

409 A more suitable framework is based on the Mutual Information (MI) which
 410 measures the dependence between two random variables X, Y and is unaffected
 411 by their marginal distributions. It is a standard approach in probability and in-
 412 formation theory (see e.g. Hlinka et al., 2014), and for discrete random variables
 413 is defined as:

$$MI(X, Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \cdot \log \left(\frac{p(x, y)}{p(x) \cdot p(y)} \right) \quad (3)$$

414 MI is zero if the two events are independent, i.e. if $p(X, Y) = p(X) \cdot p(Y)$,
 415 non-negative ($MI(X, Y) \geq 0$) and symmetric ($MI(X, Y) = MI(Y, X)$).

416 In our analysis we consider the binary variables Ψ_i at the locations i which
 417 state whether the precipitation x_i is above or below the threshold x_{0i} , i.e. $\psi_i = 1$
 418 if $x_i \geq x_{0i}$ and $\psi_i = 0$ if $x_i < x_{0i}$. Following the definition above we then
 419 calculate for each pair of locations i, j the MI for these binary variables

$$MI_{i,j} = MI(\Psi_i, \Psi_j) = \sum_{\psi_i \in [0,1]} \sum_{\psi_j \in [0,1]} (p(\psi_i, \psi_j) \cdot \log \left(\frac{p(\psi_i, \psi_j)}{p(\psi_i) \cdot p(\psi_j)} \right)) \quad (4)$$

420 We calculate MI for the dry-wet threshold $x_{0i} = 1\text{mm/d}$ as well as for a
 421 high precipitation threshold defined as the 90th percentile ($P90_i$) of the observed
 422 daily precipitation (including dry days) at each station, i.e. $x_{0i} = P90_i$.

423 Following the methodology for correlograms (see section 2.2.1), we plot each
 424 MI_{ij} against the distance of the locations i, j and fit a degree-2 loess curve to
 425 the resulting plots. We then define MI thresholds for calculating the MI lengths
 426 (MILs) for observations and for the different downscaling methods. For the dry-
 427 wet binary variable based on $x_{0i} = 1\text{mm/d}$ we use MI thresholds that depend
 428 on the experiment and season in order to obtain observed MILs that are similar
 429 (within a few kilometers) to the observed CLs, which makes it easier to assess
 430 whether MI yields information about the methods that is not already included in
 431 the CLs. The respective values are given in Table 6. For the high precipitation
 432 threshold $x_{0i} = P90_i$ we use a constant MI threshold of 0.1. Analogous to
 433 the correlation analysis MIL biases are calculated for the different downscaling
 434 methods, seasons and experiments by subtracting the respective observed MIL.

435 2.2.4 Regionalisation

436 Note that in this study, we apply the term regionalisation in the sense of spatial
 437 clustering, i.e. in the sense of finding regions with common variability. In order
 438 to achieve a regionalisation of the station data, orthogonally rotated (Varimax
 439 criterion, S-mode) principal component analysis (RPCA, e.g. Richman, 1986;

Experiments	<i>annual</i>	<i>DJF</i>	<i>JJA</i>	<i>MAM</i>	<i>SON</i>
1a, 1aRCM	0.18	0.14	0.20	0.18	0.20
1c	0.24	0.22	0.24	0.22	0.24

Table 6: MI thresholds used to calculate the MI lengths for the precipitation occurrence (1 mm threshold in the European experiments (1a and 1a-RCM) and the German experiment (1c).

440 Hannachi et al., 2007) is applied separately for each season to the correlation
 441 matrices calculated from detrended monthly timeseries.

442 The decision on the number of PCs to be rotated is based on the criterion
 443 that each retained PC has to be representative for at least one input variable,
 444 following Philipp et al. (2007). A rotated PC is considered representative for a
 445 given station if the loading of this PC at this station is larger than the loadings
 446 of the other PCs at this station by at least one standard deviation of all loadings
 447 at this station; additionally, this loading has to be statistically significant at the
 448 5% level. Each station is assigned to the region (as defined by RPCA) for which
 449 it has the highest PC loading.

450 The number of PCs is determined from observations. Then the same num-
 451 ber of PCs is used for the PCAs of the output from the downscaling methods.
 452 Following a standard approach the observed and the downscaled groupings are
 453 compared using the Adjusted Rand Index (ARI, Hubert and Arabie, 1985; San-
 454 tos and Embrechts, 2009). The ARI is based on how pairs of objects, which in
 455 our case are pairs of locations, are classified as being either in the same or in
 456 different groups, which in our case are homogeneous regions. When comparing
 457 two classifications U and V there are four options for each pair and we denote
 458 the number of pairs for each option as:

459 a number of pairs that are in the same group in both classifications

460 b number of pairs that are in the same group in U and in different groups
 461 in V

462 c number of pairs that are in the same group in V and in different groups
 463 in U

464 d number of pairs that are in different groups in U and in different groups
 465 in V

466 With these definitions, and n being the number of objects, the ARI can be
 467 expressed as

$$ARI = \frac{\binom{n}{2}(a+d) - [(a+b)(a+c) + (c+d)(b+d)]}{\binom{n}{2}^2 - [(a+b)(a+c) + (c+d)(b+d)]}. \quad (5)$$

468 Its value increases with the agreement of the two classifications; 0 indicates no
 469 agreement and the maximum is 1 for identical classifications.

470 As already mentioned in Sec. 2.2.2 the singular value decomposition routine
 471 used for PCA cannot handle missing values, and therefore the regionalisation
 472 could not be calculated for a few methods.

473 3 Results

474 3.1 Example situation

475 Before we present the results of the statistical analyses we give an example
476 for observed and downscaled precipitation on a specific day and for a few se-
477 lected methods to illustrate the different characteristics of downscaling methods
478 (Fig. 2). We chose 15. August 1998, because on this day there was frontal pre-
479 cipitation (over parts of the Scandinavia and the Baltic) as well as convective
480 precipitation (over the Iberian Peninsula and parts of northern Italy). The dis-
481 tinction is based on the analysis of pressure charts and vertical temperature
482 profiles (not shown).

483 The precipitation observations show low to medium values at most stations
484 in Northern Spain and at one station in northern Italy, while the values in Scan-
485 dinavia and the Baltic are medium to high. The ERA Interim reanalysis partly
486 underestimates the amplitudes, shows a continuous rain band south of the Alps
487 whereas only one station has recorded rainfall in this region, and does not sim-
488 ulate the convective precipitation in central Iberia. In comparison the RACMO
489 regional model simulates the intensities in some regions better, for instance over
490 Iberia and Scandinavia, but shows the well-known drizzle effect with light pre-
491 cipitation over large areas, as well as an unrealistic rain band north of the Alps
492 and over parts of Germany and France. We note that satellite pictures showed
493 convection over Germany, which however was not associated with precipitation.
494 As expected the two quantile mapping methods EQMs (empirical) and RATY
495 (parametric) inherit the partly unrealistic spatial structure from RACMO but
496 change the specific values, with the EQMs intensities being in general closer to
497 the observations than those from RATY.

498 The ANALOG-ANOM method captures well the fact that the convective
499 precipitation only occurs at some locations and that the frontal precipitation is
500 more homogeneous in space. The individual locations at which the convective
501 precipitation occurs are partly different to the observations, which is an expected
502 consequence of the stochastic nature of occurrence of convection. The values for
503 the convective precipitation are close to the observed ones, whereas the intensity
504 of the frontal precipitation is underestimated.

505 The MLR-RAN method (PP, multiple linear regression using large-scale pre-
506 dictors) unrealistically yields precipitation at all locations with the exception of
507 some stations close to the eastern boundary of the analysis domain. For the sta-
508 tions where precipitation was observed the intensities are roughly in the right
509 range. For the WT-WG method (weather generator conditioned on weather
510 types) one can either plot individual realisations or the average over a simu-
511 lated ensemble (100 realisations in this case). The individual realisations (not
512 shown) have a much too low spatial coherency. This indicates that the random
513 variability component, which is sampled individually at each location, is large
514 compared to the fraction of variability that is conditional on the weather types.
515 Here we show the conditioned component, i.e. the averaged values, which is
516 as expected, too smooth, with precipitation occurring almost everywhere and
517 values at the locations with observed precipitation being often too low.

518 In summary, the examples suggest that the methods that either inherit the
519 spatial structure from an RCM (EQMs and RATY) or use observed spatial struc-
520 tures (ANLOG-ANOM) yield relatively realistic spatial patterns. In contrast

521 conditioning precipitation at single sites on large-scale predictors (MLR-RAN,
 522 WT-WG) leads to fields that are too smooth when only the conditioned com-
 523 ponent is considered (MLR-RAN, averaged WT-WG), or not smooth enough
 524 when the stochastic component is added (individual realisations of WT-WG).

525 3.2 Correlations

526 Selected examples of pairwise cross-correlation matrices for winter (DJF) are
 527 displayed in Figs. 3a and 3b for precipitation and maximum temperature res-
 528 pectively. The 86 European stations (Fig. 1) are arranged so that station pairs
 529 with a small distance are near to the diagonal while distant pairs are near the
 530 upper-left corner. The geographic distances (measured along a great circle) are
 531 shown in the upper triangle in the first matrix of each panel, while the observed
 532 correlations are shown in the lower triangle.

533 In general, all methods are able to reproduce to some extent the correlation
 534 structure of both temperature and precipitation, with the exception of WT-WG.
 535 The WT-WG correlations shown are the average of the correlations for individ-
 536 ual realisations (in contrast to Fig. 2 where correlations for ensemble-averaged
 537 values are shown), and despite the conditioning of the weather generator on
 538 weather types it yields almost uncorrelated values for all stations, regardless of
 539 their distance. This is explained by the weak conditioning imposed by the only
 540 predictor (SLP) used in this method, which explains only a very small frac-
 541 tion of the variance and results in an almost purely stochastic method (see also
 542 Gutiérrez et al. (2018)). The correlations of raw ERA-Interim and RACMO
 543 output are both in good overall agreement with the observations. However, the
 544 results for the different methods differ in detail. For instance, MLR-RAN sys-
 545 tematically yields too high positive and negative correlations for distant station
 546 pairs, while EQMs and in particular ANALOG-ANOM reproduce most aspects
 547 of the structure well. The latter has the highest CMD value for precipitation
 548 (0.988) and maximum temperature (0.992).

549 We now investigate the dependency of correlations on distance more system-
 550 atically by comparing correlograms and CL values. The former are shown for
 551 some example methods in Fig. 4 for the European station network (experiments
 552 1a and 1aRCM) and in Fig. 5 for the high-density German network (experiment
 553 1c). In addition to the actual correlations these figures include the fitted curves
 554 and the CLs (vertical lines). As expected the observed correlations (upper-left
 555 panels) decline with distance and for the European dataset level off around zero.
 556 The fact that the correlations show approximately an exponential decrease in
 557 Fig. 4 but a more linear decrease in Fig. 5 is due to the different size of the
 558 analysis domains. In experiment 1c there are some missing CL values for tem-
 559 peratures, because due to the small analysis domain and the smooth topography
 560 the temperature records are highly correlated for all station pairs and in some
 561 cases the fitted line is therefore above the corresponding correlation threshold
 562 (0.65, Table 3) for all distances. In contrast precipitation has a higher degree
 563 of spatial heterogeneity and CLs are obtained in all cases.

564 For the European data (Fig. 4) ERA-Interim tends to slightly overestimate
 565 the correlations in both seasons and reproduces the observed slight difference be-
 566 tween the seasons. RACMO has values closer to reality, but does not capture the
 567 observed seasonal difference. Both MOS methods (EQMs-R and Ratyetal-M8-
 568 R) further reduce the correlations compared to the raw RCM but to a different

569 extent, and the lack of a seasonal difference remains. As expected, the ana-
 570 log method (ANALOG-ANOM), which selects an entire analog field reproduces
 571 the observed correlations. The PP example method (MLR-RAN), which uses
 572 large-scale predictors, overestimates correlations. As already shown in Fig. 3a
 573 the weather generator conditioned on weather types (WT-WG) strongly un-
 574 derestimates correlations when individual realisations are considered. For the
 575 ensemble average (dashed lines) the correlations are too high in winter and still
 576 substantially too low in summer. The deficiencies of this method have also been
 577 reported in Gutiérrez et al. (2018).

578 It can be seen in Fig. 5 that in Germany and on the shorter distances, which
 579 are resolved well by the high-density network, the observed seasonal differences
 580 are larger than in the European case, with higher correlations in winter. All ex-
 581 ample methods do now also show a seasonal difference. As in the European case
 582 ERAINT overestimates correlations. The MOS-corrected ERA-Interim precip-
 583 itation (EQM-R) leads to fairly realistic correlations, as does one of the PP
 584 methods (DSCLIM-D), while the other ones either overestimate (GLM-DET)
 585 or underestimate (GLM-BN-DET) correlations. As explained in section 2.1,
 586 the latter is an extension of the former, explicitly including a model for spatial
 587 dependence (based on probabilistic networks).

588 We now look at the full set of methods with respect to the precipitation CL
 589 bias for the European (Fig. 6) and German datasets (Fig. 7). In Fig. 6 ERAINT
 590 has a positive CL bias, which gets reduced when the reanalysis is dynamically
 591 downscaled with RACMO, as already seen in the previous figures. Most deter-
 592 ministic MOS methods do reduce the bias both in the reanalysis-driven (*-E)
 593 and RACMO-driven (*-R) case, with the former still having higher CLs than the
 594 latter, as for the raw numerical models. Many MOS methods that are based on
 595 quantile mapping have very low CL biases, while some of the scaling approaches
 596 (e.g. Ratyetal-M7) have slightly higher biases. Consistent with the previous
 597 plots, the stochastic methods (MOS-GLM, VGLMGAMMA) have substantial
 598 negative CL biases for the individual realisations. The bias for the ensemble
 599 mean is positive for MOS-GLM, while it is negative for VGLMGAMMA, sug-
 600 gesting that for the latter the distributions are not constrained closely enough
 601 by the predictands.

602 The PP methods in Fig. 6 show a wide range of positive and negative bi-
 603 ases. Positive biases occur for regression methods with large-scale predictors
 604 (MLR-RAN, MLR-RSN, MLR-ASI, GLM-DET) because the predictors for dif-
 605 ferent stations are similar (e.g. PCs from ERA-Interim fields). The FIC01P
 606 method, which is a combination of an analog method and postprocessing us-
 607 ing a transfer function, has also a positive bias. In contrast, negative biases
 608 are visible for methods that use local predictors, e.g. information taken from
 609 the gridcell covering the target station, for instance some of the linear mod-
 610 els (GLM, GLM-WT, GLM-P) and the 'multi-objective genetic programming
 611 method' (MO-WT). The ANALOG method, which is based on regional-scale
 612 predictors shows a negative CL bias. Individual realisations of some stochastic
 613 methods (ANALOG-M, ANALOG-SP, GLM-P) have also negative biases. Bi-
 614 ases close to zero are achieved with one analog method (ANALOG-ANOM) and
 615 a regression method with noise added (MLR-ASW).

616 When the CL biases on shorter distances are considered (Fig. 7) the raw
 617 ERA-Interim precipitation shows again a positive bias, while biases close to zero
 618 are obtained for MOS methods based on quantile mapping. For the PP methods

619 the positive biases of regression methods using large-scale predictors and the
620 negative bias for those using local predictors remain. The ANALOG method is
621 now almost bias-free, in contrast to the European case. The reason is that the
622 predictors are neither global, nor completely local, but based on the division
623 of the whole domain in a number of sub-domains with each containing several
624 stations. The selection of analog dates is common for all stations within a sub-
625 domain, thus guaranteeing the spatial consistency within sub-domains, whereas
626 different dates can be chosen for different sub-domains. As Germany lies within
627 one sub-domain and Europe covers several subdomains the CL bias is close to
628 zero for experiment 1c (sampling effects remain) and negative for experiment
629 1a. The second method that is bias-free is a hybrid method (DSCLIM-D) which
630 combines a weather type based transfer function and an analog approach.

631 For the European dataset we also consider the CL bias for minimum and
632 maximum temperature (Figs. 8 and 9). As temperature fields are smoother than
633 precipitation fields, we use a correlation threshold of 0.5 rather than 0.35, which
634 was used for the European precipitation data. The results for minimum and
635 maximum temperatures are very similar. The MOS results are fundamentally
636 different from the precipitation case. While for precipitation many MOS meth-
637 ods did reduce the CL bias relative to the raw models (both for ERAINT and
638 RACMO), for temperature there is for almost all MOS methods no reduction of
639 the positive model bias. The reason might be that precipitation is an intermit-
640 tent process for which debiasing the marginal distribution affects correlations
641 more strongly than for the continuous temperature timeseries. The high biases
642 for CDFt-E and MOS-REG-R need further investigation. The CDFt method
643 was also found to behave differently to other MOS methods with respect to
644 the temporal correlation between predictions and observations (Gutiérrez et al.,
645 2018), trends (Maraun et al., 2018) and extreme events (Hertig et al., 2018).
646 We note that this method is different from the other MOS techniques in the
647 sense that it also uses the predictand distribution in the validation period (see
648 Gutiérrez et al. (2018), Appendix A.1 for the full method description), which
649 may lead to a high sampling variability in our experimental setup. The CDFt
650 data passed our standard quality test, but the correlation vs. distance plots for
651 CDFt for maximum and minimum temperatures and experiment 1a showed an
652 unusual behaviour with no clear link between correlations and distance, and
653 thus a technical error for downscaled temperatures using CDFt-E cannot be
654 ruled out.

655 As for precipitation the PP methods show again in general higher biases than
656 the MOS methods, and some analog methods perform well, whereas others do
657 not. A noticeable difference is the smaller number of methods with negative
658 biases for temperature. Although the set of methods is not identical, there are
659 some methods used for both predictor variables that have large negative biases
660 for precipitation but small biases for temperature, namely ANALOG-SP and
661 MO-GP. A potential reason is that for those methods the predictors constrain
662 temperature better than precipitation.

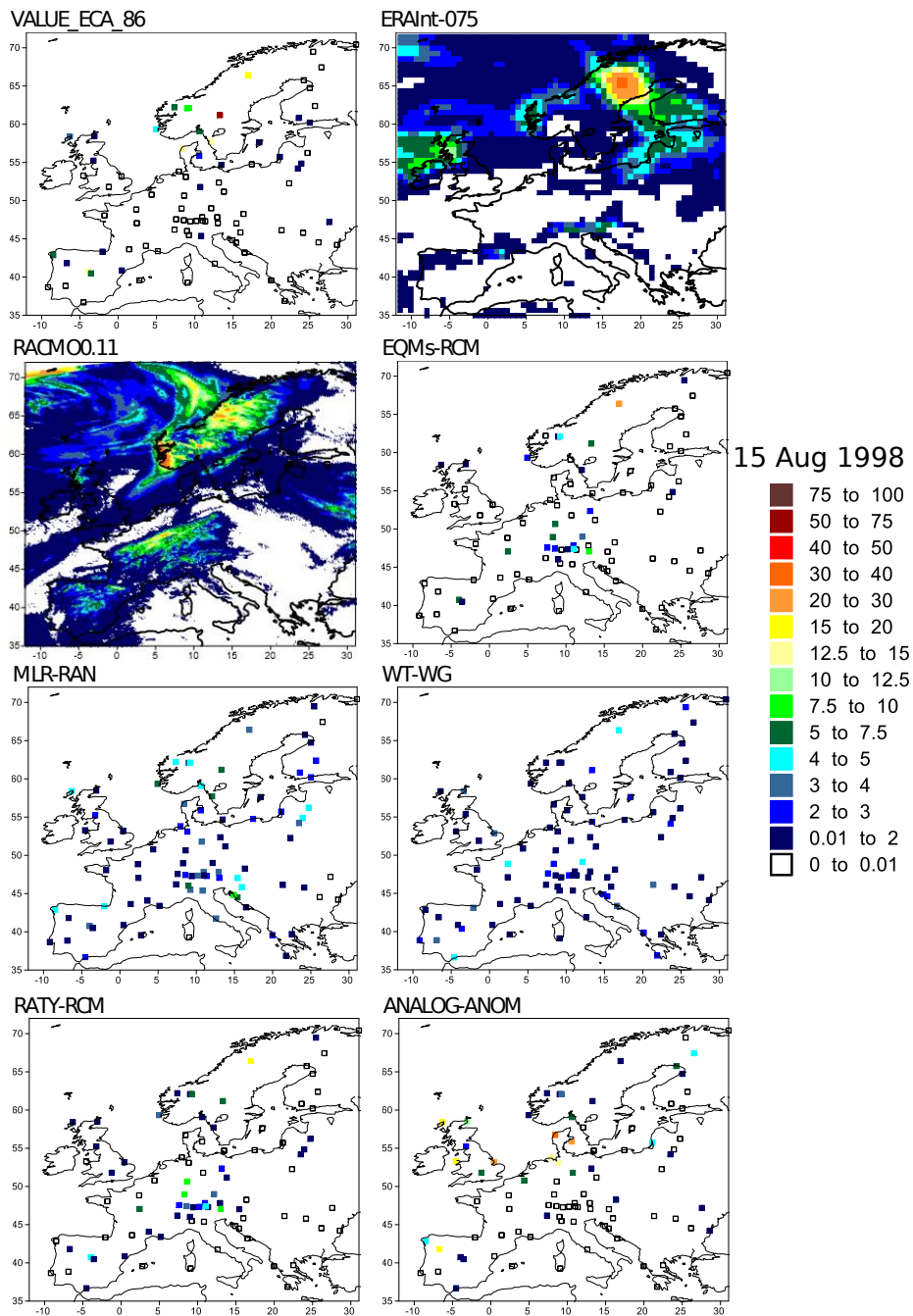
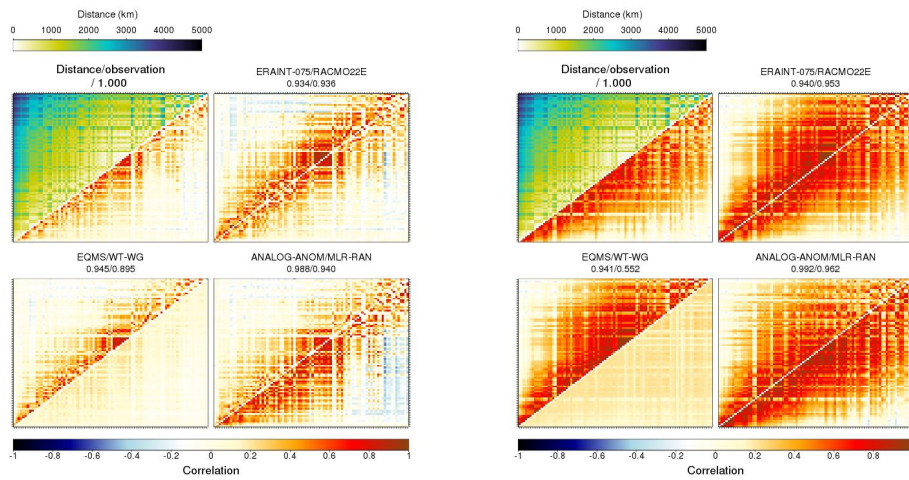


Figure 2: Observed (VALUE-ECA-86-v2, top-left panel) and downscaled precipitation on 15. August 1998 (mm/d). The second and third panels (from top to bottom, and left to right) show the 24h accumulated precipitation from the ERA-Interim reanalysis (ERAInt-075 panel) and from the RACMO RCM (0.11 degree horizontal resolution, RACMO 0.11 panel) driven by ERA-Interim. The downscaling methods are labelled by their codes (Table 1), with the “-RCM” suffix indicating MOS methods used in experiment 1a-RCM.



(a) Daily DJF precipitation (Spearman's ρ correlation coefficient)

(b) Daily DJF maximum temperature (Pearson's r correlation coefficient)

Figure 3: Pairwise cross-correlation matrices for winter for the 86 locations of the VALUE-ECA-86-v2 dataset. In each panel, the first matrix represents the geographic distances between pairs of stations (above the diagonal) and the correlations of the observations (below the diagonal). The remaining matrices display the correlations for two different methods indicated by the panel titles with the values for the first (second) method given above (below) the diagonal. The number under the method names is one minus the correlation matrix distance between the method and the observation correlation matrices.

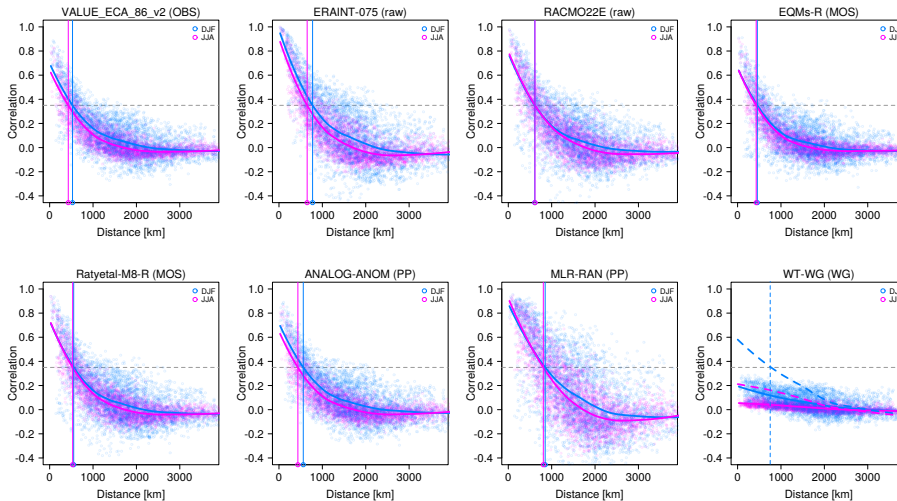


Figure 4: Correlograms for daily precipitation for JJA and DJF showing correlations of the timeseries for each pair of stations against their geographical distances (European experiment, expla). For the stochastic WT-WG method the fitted curves of the *averaged* option and the corresponding CL value are indicated by dashed lines (individual values are omitted for clarity).

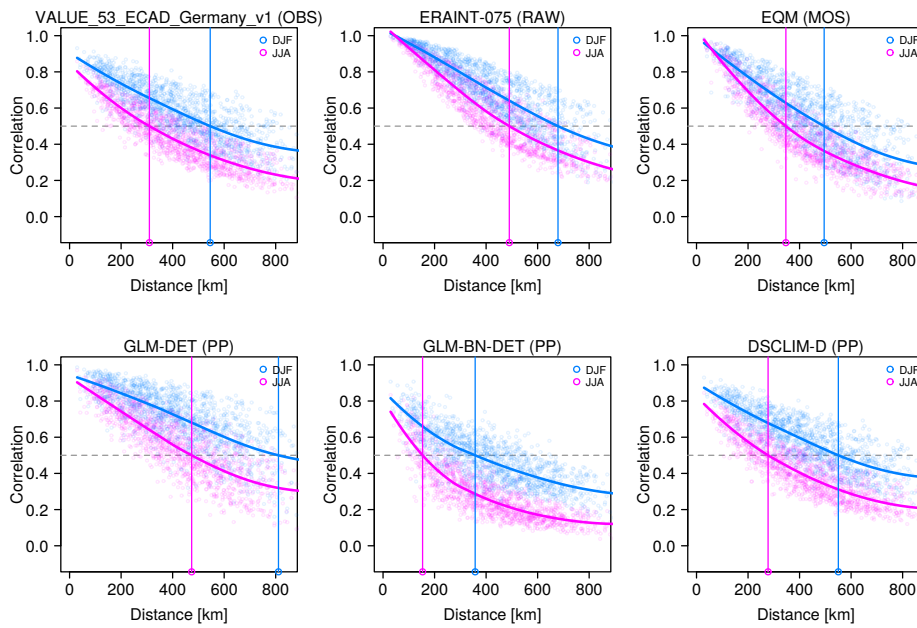


Figure 5: Same as Fig. 4 but for selected methods used in the German experiment (exp1c). The correlations for the reference observations (VALUE-ECA-53-Germany-spatial-v1) are shown in the upper left panel.

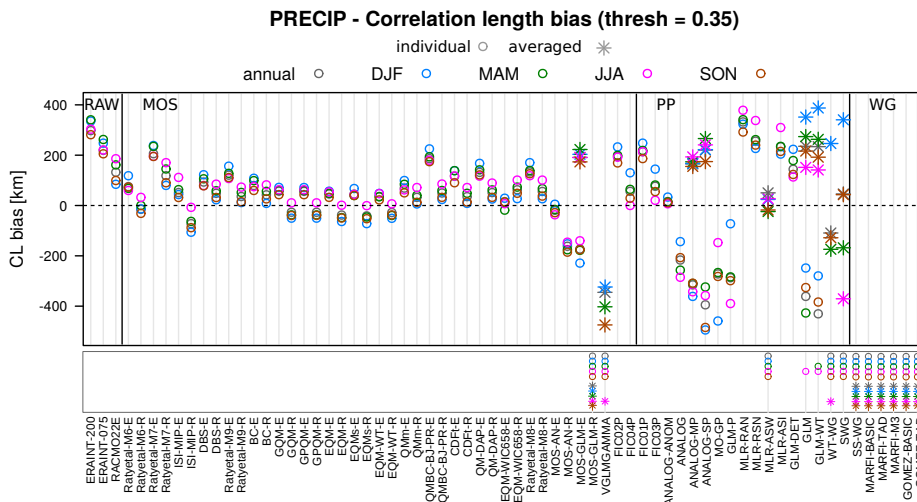


Figure 6: Correlation length (CL) biases for daily precipitation from the downscaling methods tested in experiments 1a (suffix $-E$ for MOS methods) and 1a-RCM (suffix $-R$) with respect to the reference values based on the VALUE-ECA-86-v2 dataset (Table 4). For the stochastic methods, the results of both the member-averaged (asterisks) and individual (circles) approaches are shown. The box in the lower part of the figure shows the seasons/approaches for which the CL cannot be calculated due to very low correlations.

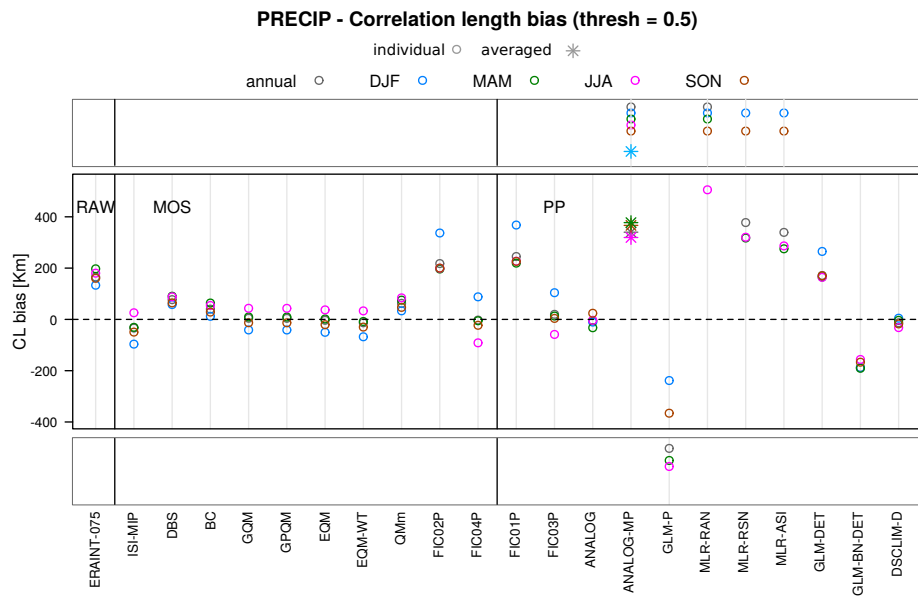


Figure 7: Same as Fig. 6 but for the German experiment (exp1c). The boxes in the lower/upper part of the figure show the seasons and approaches for which CL distance cannot be calculated. Upper box: the fitted correlogram line is entirely above the threshold. Lower box: the fitted correlogram line is entirely below the threshold.

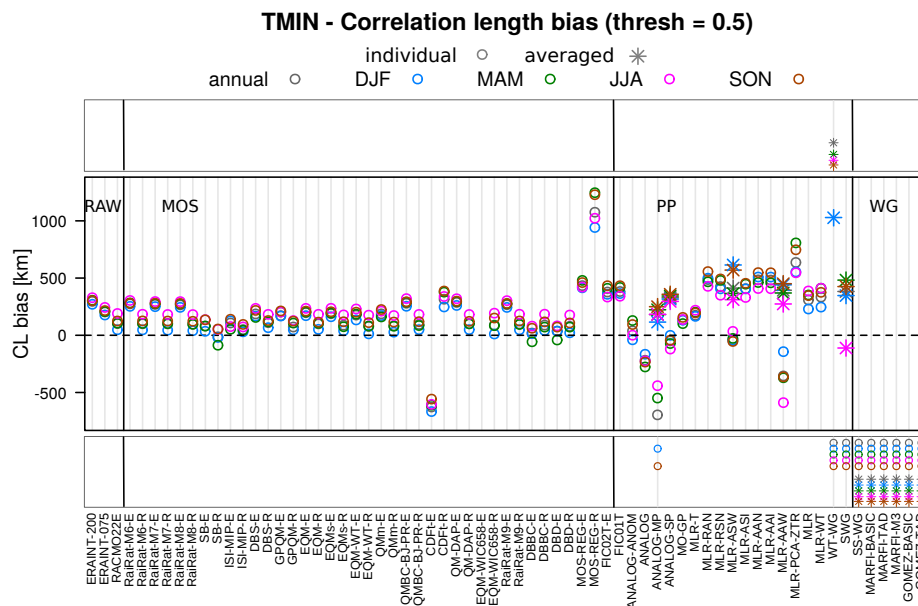


Figure 8: Same as Fig. 6 but for minimum temperature.

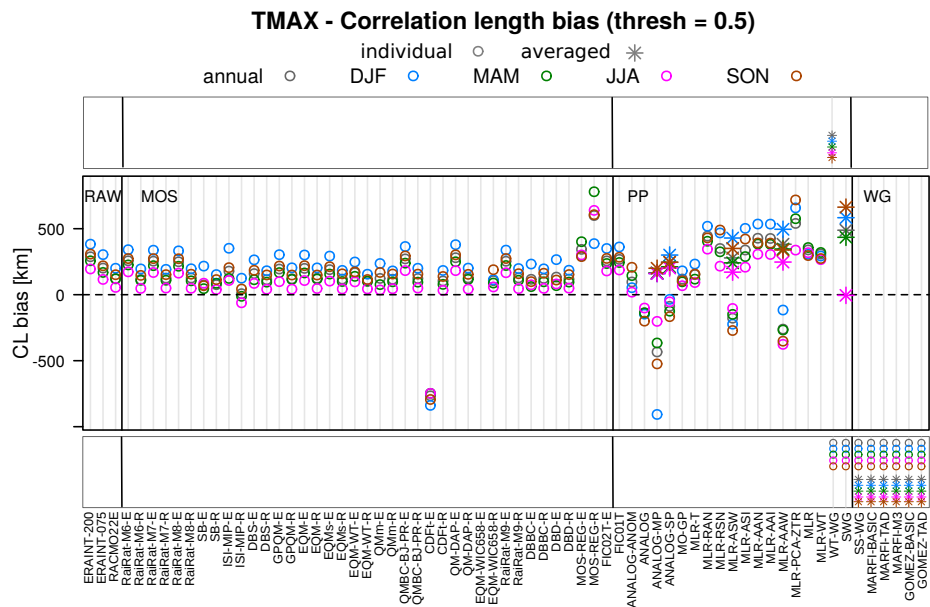


Figure 9: Same as Fig. 8, but for maximum temperature.

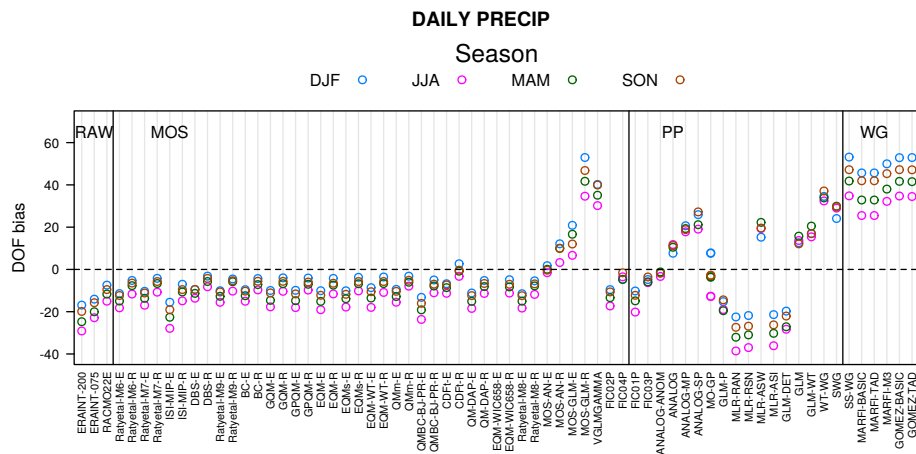


Figure 10: Bias of the spatial degrees of freedom (DOF) for daily precipitation from the methods included in the European experiments (expla and expla-RCM).

663 **3.3 Spatial degrees of freedom**

664 The DOF biases for precipitation, which express differences in the dimensionality of the fields, are shown in Fig. 10. Almost all MOS methods have a negative bias and thus underestimate complexity. The underestimation is strongest in summer, where convective, and thus small-scale, precipitation is more important than in the other seasons. Compared to the raw model results, most MOS methods reduce the absolute bias. The exception are some of the stochastic methods (MOS-GLM, VGLMGAMMA), which strongly overestimate complexity. The MLR-based PP methods also underestimate complexity, whereas some of the analog methods have a small bias and others overestimate it. The weather generators show a strong overestimation.

674 The DOF biases for temperature are shown in Fig. 11. For almost all downscaling methods they are substantially smaller than for precipitation, with many MOS and some PP methods leading to biases smaller than 2. The exception are some WG methods (SS-WG, GOMEZ-BASIC, GOMEZ-TAD), which show biases of up to 40. During summer and autumn the DOF biases for minimum temperature are larger than those for maximum temperature. In contrast to the precipitation case the biases for the MOS-corrected models are very similar to those of the raw models.

682 Most methods with a positive (negative) CL bias, i.e. those for which correlations drop too slowly (too quickly), have a negative (positive) DOF bias. One clear exception is CDFt-E for temperature, which is in line with other MOS methods with respect to the underestimation of the DOFs, but as mentioned in Sec. 3.2 has a large negative CL bias, which may be due to technical errors. We note that reordering the stations would not affect the DOFs, but would lead to erroneous correlograms if not taken into account when calculating the distances between station pairs. There are also some MOS methods that have a slightly positive CL bias despite their negative DOF bias.

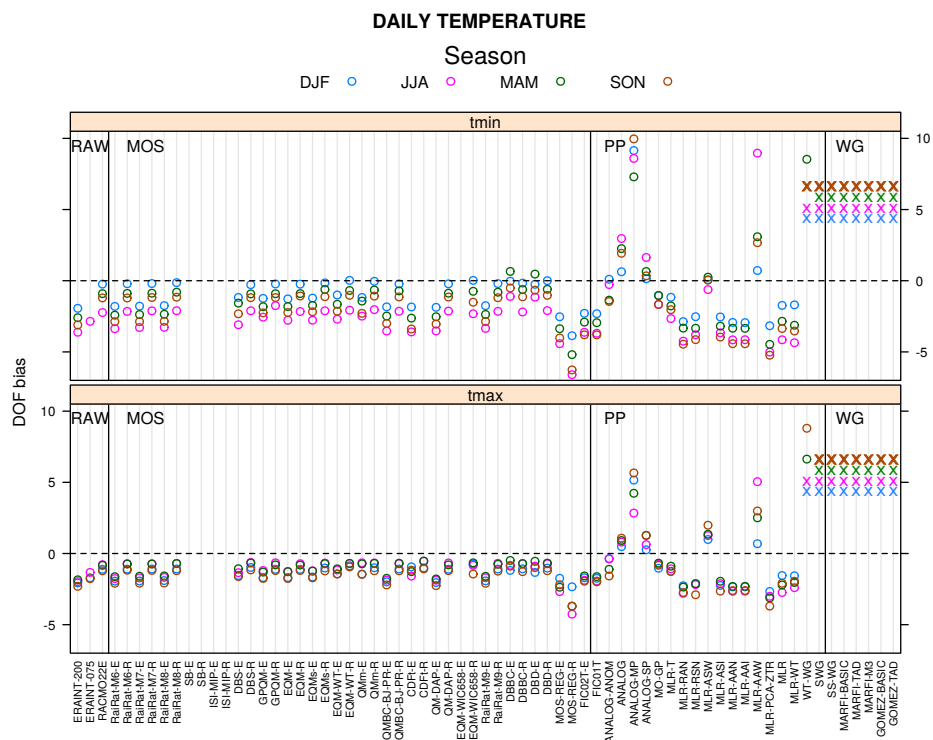


Figure 11: Same as Fig. 10, but for daily temperature. The methods marked with a cross (×, coloured according to the season) are out of range with positive bias of more than 10 degrees of freedom. The methods without results are those having missing values in the covariance matrix (see Sect. 2.2.2).

691 3.4 Joint threshold exceedances

692 The methodology for the joint threshold exceedances analysis is very similar
693 to that for correlation (see Sec. 2.2.3), and we therefore do not show the MI
694 matrices and MI vs. distance diagrams. The characteristic MI lengths for the
695 reference observations exceeding the wet day threshold are presented in Tab. 7
696 and for exceeding the local 90th percentile in Table 8. As in the case of the
697 correlograms, lower MIL values indicate a faster loss of mutual dependence as
698 a function of distance, while higher MIL values indicate a stronger dependence
699 between stations. For both thresholds there is a marked seasonal dependence,
700 with the minimum in summer and the maximum in winter. For the 90th per-
701 centile autumn values are also high. The MILs obtained from the European and
702 the German observational datasets were similar (Table 7).

703 The high-density German dataset is better suited than the European dataset
704 for calculating MILs for both thresholds, as it has a larger number of station
705 pairs within the distance ranges relevant for calculating the MILs for both
706 thresholds, and thus provides more robust results. We therefore restrict the
707 MIL analysis to experiment 1c. This has the additional advantage that we
708 avoid a potential loss of robustness in the summer results arising from locations
709 with no precipitation for the whole season, which may occur in some parts of
710 Southern Europe. The biases for the wet day threshold with respect to the ob-
711 served reference values are shown in Fig. 12 and for the 90th percentile threshold
712 in Fig. 13.

713 For the wet day threshold all MOS methods slightly overestimate the depen-
714 dence. The exceptions are FIC02P, which strongly overestimates it, and FIC04P,
715 which in most seasons slightly underestimates it. All MOS methods but FIC02P
716 reduce the bias compared to the raw reanalysis data. Among the PP methods
717 ANALOG and DSCLIM-D (which contains an analog step) are bias-free apart
718 from sampling effects, and the individual realisations of ANALOG-MP has also
719 a very low bias. The MLR methods overestimate the dependence, whereas
720 GLM-P strongly underestimate it.

721 The different downscaling methods perform similarly with respect to the
722 MIL biases for the wet day threshold and to the CL biases (Fig. 7). Both show
723 a bias reduction by most MOS methods, and the same sign and relative size of
724 the bias for both quantities. Too strong (weak) correlations of the timeseries
725 are thus associated with too high (low) dependences of the occurrence of wet or
726 dry days.

727 The overall picture is different for the 90th percentile threshold. Almost all
728 MOS methods show the same overestimation of dependence as the raw reanalysis
729 data. In the PP group the analog methods and GLM-BN-DET and DSCLIM-
730 D have very low biases, whereas the MLR methods very strongly overestimate
731 dependences for heavy precipitation.

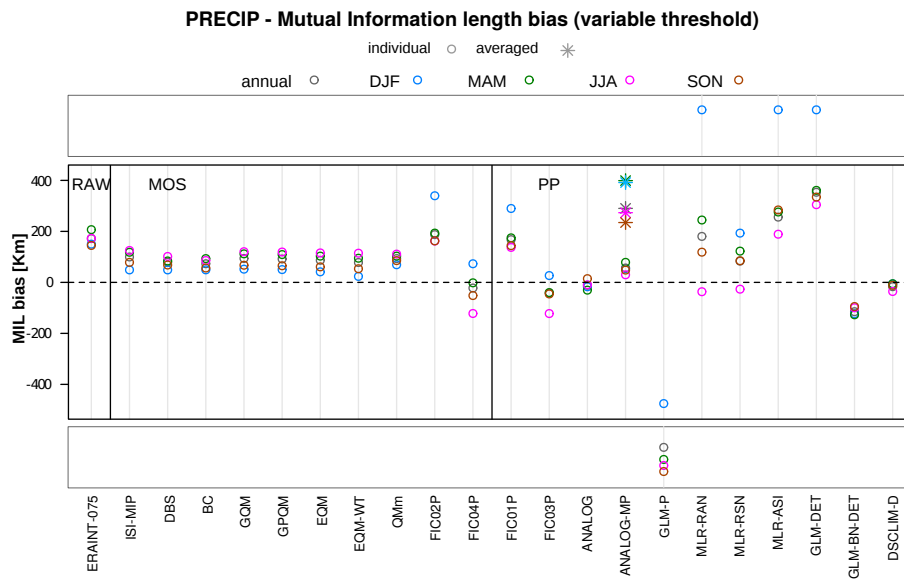


Figure 12: Mutual Information Length (MIL) biases for the exceedance of the wet day precipitation threshold, obtained from experiment 1c with respect to the values from the reference observations (VALUE-53-ECAD-Germany-v1 dataset, Table 7). Values in the boxes in the upper and lower part of the figure indicate methods for which the MI Length value cannot be calculated due to the MI values being too high or too low (as for correlations in Fig. 7).

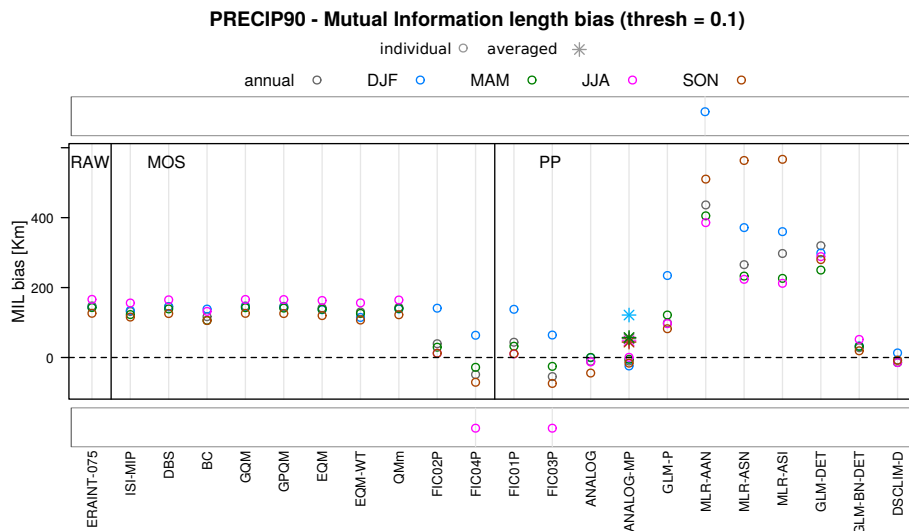


Figure 13: Same as Fig. 12 but for the exceedance of the 90th percentile of daily precipitation obtained from experiment 1c.

Experiments	<i>annual</i>	<i>DJF</i>	<i>JJA</i>	<i>MAM</i>	<i>SON</i>
1a, 1aRCM	359	554	216	324	340
1c	338	528	256	360	345

Table 7: Mutual Information Length (MIL) values (in km) calculated for exceedance of the wet day threshold of daily precipitation in the reference station datasets (VALUE-ECA-86-v2 for experiments 1a and 1a-RCM and VALUE-ECA-53-Germany-spatial-v1 for experiment 1c), using the MI thresholds displayed in Tab. 6.

<i>annual</i>	<i>DJF</i>	<i>JJA</i>	<i>MAM</i>	<i>SON</i>
191	284	109	183	234

Table 8: Mutual Information Length (MIL) values (in km) calculated for exceedance of the 90th percentile of daily precipitation in the reference station dataset of experiment 1c (VALUE-ECA-53-Germany-spatial-v1, using a fixed threshold of 0.1 for all seasons. Note that only the experiment 1c (German dataset) has been used in this case as reference (see Sec. 3.4).

732 3.5 Regionalisation

733 The number of PCs retained for rotation is shown in Table 9 along with the
 734 cumulative fraction of variance explained for the observed daily precipitation,
 735 minimum and maximum temperature at the 86 European stations. As expected
 736 a higher number of PCs is needed to explain a certain fraction of the variability
 737 of precipitation compared to temperature, as the spatial patterns of the former
 738 contain more small-scale structures. Also, more PCs are needed to represent
 739 precipitation well in summer and spring than in autumn and winter, due to
 740 the higher contribution of small-scale, convective precipitation in the former
 741 seasons, and the dominance of large-scale, stratiform precipitation in the latter.
 742 The fact that the retained PCs do not explain all the variance in the datasets
 743 is one of the potential reasons for differences between the rotated EOFs in the
 observations and the downscaling results.

Var.	<i>DJF</i>	<i>MAM</i>	<i>JJA</i>	<i>SON</i>
Precip	13 (78.3)	15 (71.4)	19 (71.4)	13 (71.8)
Tmin	6 (85.8)	6 (85.1)	6 (82.3)	6 (81.4)
Tmax	6 (87.2)	6 (87.3)	6 (86.5)	5 (82.0)

Table 9: Number of principal components retained for rotation and cumulative variance (in parentheses, %) for precipitation, minimum and maximum temperature at the 86 stations of the ECA-VALUE-86-v2 observation dataset.

744 For temperature 5-6 PCs are retained and thus 5-6 regions are identified.
 745 The regions for maximum temperature in the different seasons are shown in
 746 Fig. 14. Europe is divided roughly into northern Europe, north-western Eu-
 747 rope, south-western Europe, central and southern Europe, eastern Europe, and
 748 south-eastern Europe. The boundaries between the regions are to some extent
 749 seasonally dependent. They are also not always simply connected geographical
 750 regions, as for instance in autumn and spring one station in northern Italy is
 751 grouped together with the south-western stations, or in winter the UK, Germany
 752 and the Alpine regions contain stations associated with different rotated PCs.
 753 Similar regions are found for minimum temperature, but there are also some dif-
 754 ferences, for instance a distinct central alpine region for minimum temperature
 755 in winter (not shown).
 756

757 Fig. 15 shows the ARI for minimum and maximum temperatures, which is
 758 used as performance measure to judge the ability of the downscaling methods
 759 to capture the observed regions of similar temperature variations. It can be
 760 seen that the single-site WG based methods (GOMEZ-BASIC, GOMEZ-TAD,
 761 MARFI-BASIC, MARFI-TAD, MARFI-M3, SS-WG) are not able to reproduce
 762 the regions at all due to the generation of synthetic time series at one specific
 763 location without considering spatial relationships. WG methods that include at-
 764 mospheric covariates (WT-WG, SWG) perform somewhat better by indirectly
 765 incorporating spatial information carried by the covariates. There is no system-
 766 atic difference between MOS and PP methods. The ARI mostly lies between
 767 about 0.3 and 0.9 and varies more between seasons than between methods. The
 768 best performance is achieved for spring to autumn, whereas in winter the lowest
 769 ARI values are systematically attained. The lower performance in winter might
 770 partly be explained by region-specific phenomena (for instance inversion), which

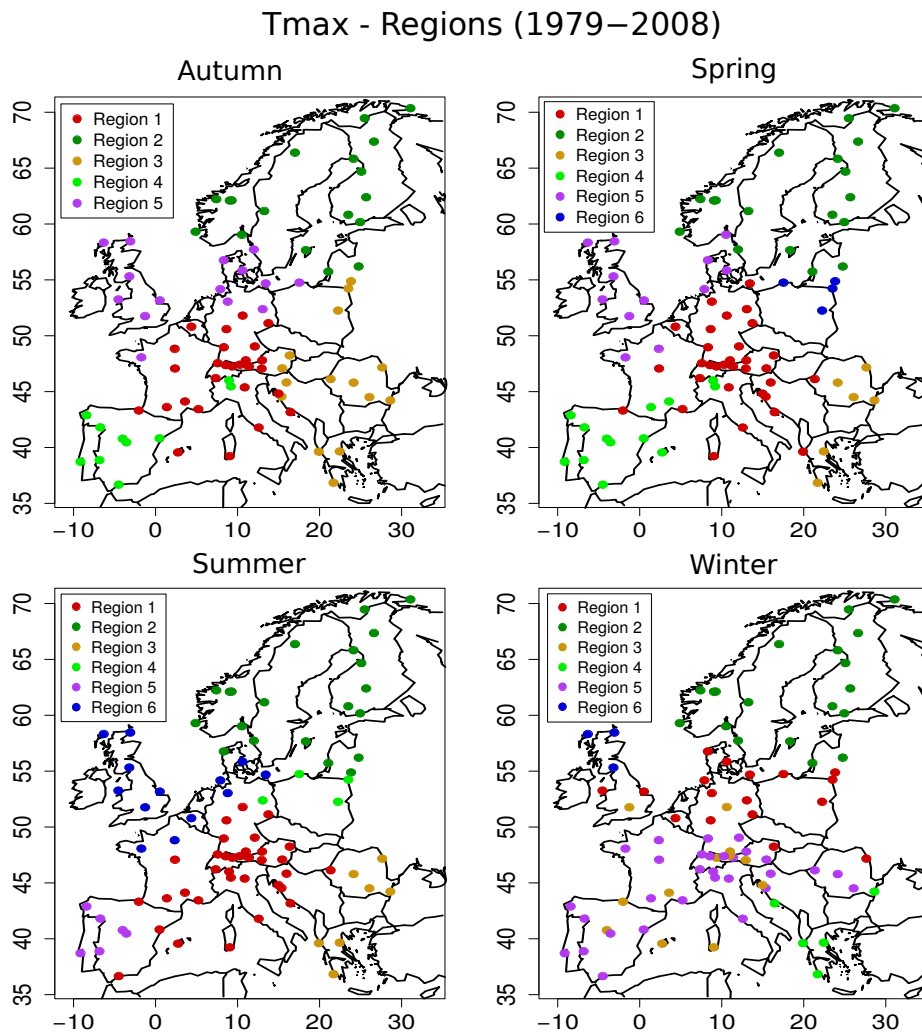


Figure 14: Regions derived from rotated PCA of seasonally detrended monthly maximum temperatures in the period 1978-2008, considering the 86 stations of the VALUE-ECA-86-v2 observational dataset.

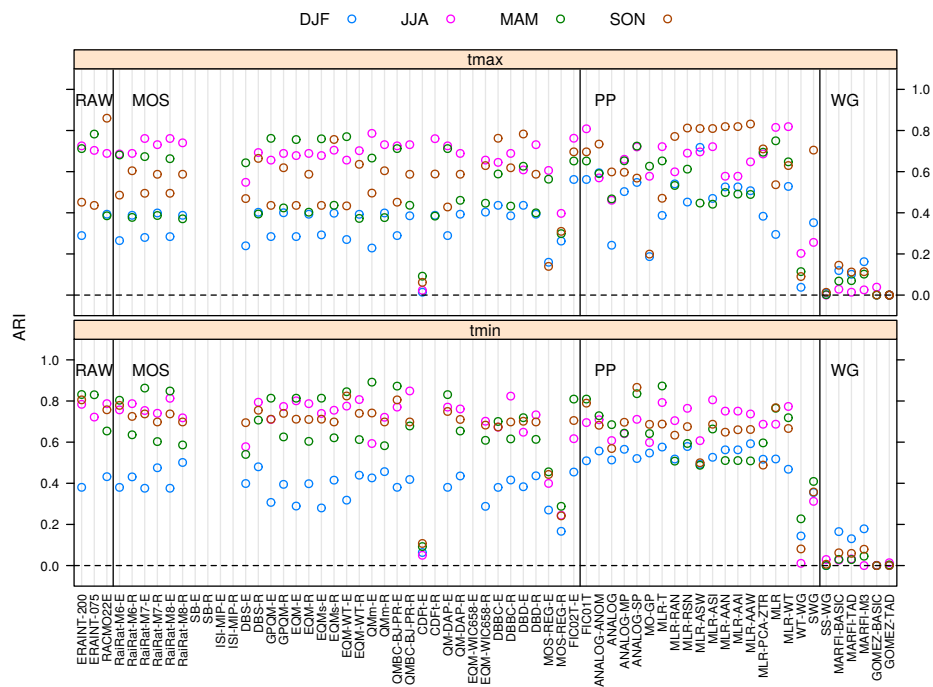
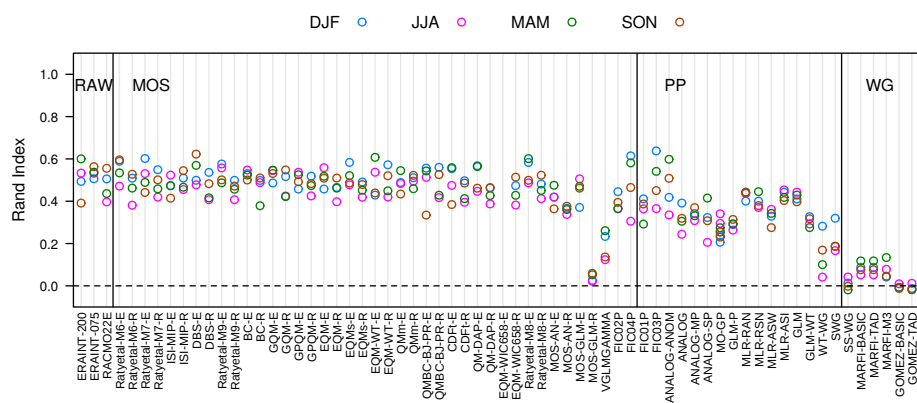


Figure 15: Adjusted Rand Index (ARI) for minimum (top) and maximum (bottom) temperatures obtained from the European experiments (exp1a and exp1a-RCM). ARI measures the agreement between the regionalisations for the observations (VALUE-ECA-86-v2 stations) and the downscaling output, ranging from 0 (no agreement) to 1 (perfect agreement). The methods without results are those having some missing values in the covariance matrix, as indicated in Section 2.2.2.

771 are not adequately captured by the downscaling methods. The ARI for ana-
 772 log methods, which by construction lead to a realistic spatial structure of the
 773 daily fields, is not higher than for many other methods. The monthly means to
 774 which the rotated PCA is applied, might be somewhat different from the true
 775 monthly means, and the questions arises to what extent the results of the ro-
 776 tated PCA describe robust statistical properties, and to what extent they might
 777 be influenced by the individual realisations. The ARI for precipitation is shown
 778 in Fig. 16 and lies between about 0.2 and 0.6, but with no seasonal structure to
 779 it. Like for temperature, WGs are not able to map the regions and no superior
 780 performance of multi-site methods arises (not shown).



781 4 Summary and conclusions

782 We have evaluated the spatial variability of the output from over 40 downscaling
783 methods for the period 1979-2008 at a European-wide network of 86 stations,
784 and at a high-resolution network of 53 stations in Germany. Predictors for the
785 PP methods and boundary conditions for the RACMO regional model have
786 been taken from the ERA-Interim reanalysis. MOS methods have been ap-
787 plied to the reanalysis as well as to the RACMO output. We have analysed
788 the dependency of correlations of daily temperature and precipitation series at
789 station pairs on the distance between the stations. For the European dataset
790 we have also investigated the complexity of the downscaled data by calculating
791 the number of independent spatial degrees of freedom. For daily precipitation
792 at the German network we have additionally evaluated the dependency of the
793 joint exceedance of the wet day threshold and of the local 90th percentile on
794 the distance between the stations. Finally we have investigated regional pat-
795 terns of European monthly precipitation and temperature obtained from rotated
796 principal component analysis.

797 The results for correlation lengths and degrees of freedom based on the Eu-
798 ropean network are summarised in Fig. 17. Findings related to joint threshold
799 exceedances are not included in the figure because they are based on the German
800 predictand data and a different set of methods. Results from the regionalisa-
801 tion are not included because they are derived from monthly rather than daily
802 data. The figure shows the relative bias calculated as the ratio of the bias and
803 the observed value for the correlation lengths or the degrees of freedom. This
804 normalisation makes it easier to compare the values for differen seasons, and
805 for correlation lengths and degrees of freedom. For the bias in the degrees of
806 freedom we have swapped the sign because a bias in correlation lengths is usu-
807 ally associated with a bias of the opposite sign in the degrees of freedom. The
808 summary figure and the detailed results presented earlier show that there is a
809 very large spread in how well the different downscaling methods represent the
810 characteristics of the observations, ranging from close to reality to very unreal-
811 istic.

812 For all three predictand variables the raw models have positive biases in
813 correlation length and negative biases in the number of degrees of freedom. The
814 biases for the RACMO model are smaller than those for the reanalysis, which
815 demonstrates the benefit of the explicit representation of smaller spatial scales.
816 It is likely that these biases are not fully due to model deficiencies because the
817 spatial scales of the data are different. Observations averaged over the gridcells
818 can have higher correlations between two locations than local values, and the
819 number of degrees of freedom of spatial averages can be lower. Likewise the
820 dependence of the exceedance of thresholds at different locations, for which the
821 models showed a positive bias, might be higher for area means than for local
822 values. Nevertheless the biases represent actual errors if the gridcell values are
823 used as direct estimates for local values.

824 As can be seen in Fig. 17 most MOS methods substantially reduce the posi-
825 tive biases in correlation length for precipitation, whereas there is no clear
826 improvement for temperature. This difference might be due to the fact that
827 precipitation is an intermittent process with many zero values, for which cor-
828 recting the simulated marginal distribution affects correlations and threshold
829 exceedances more strongly than for the continuous temperature timeseries. The

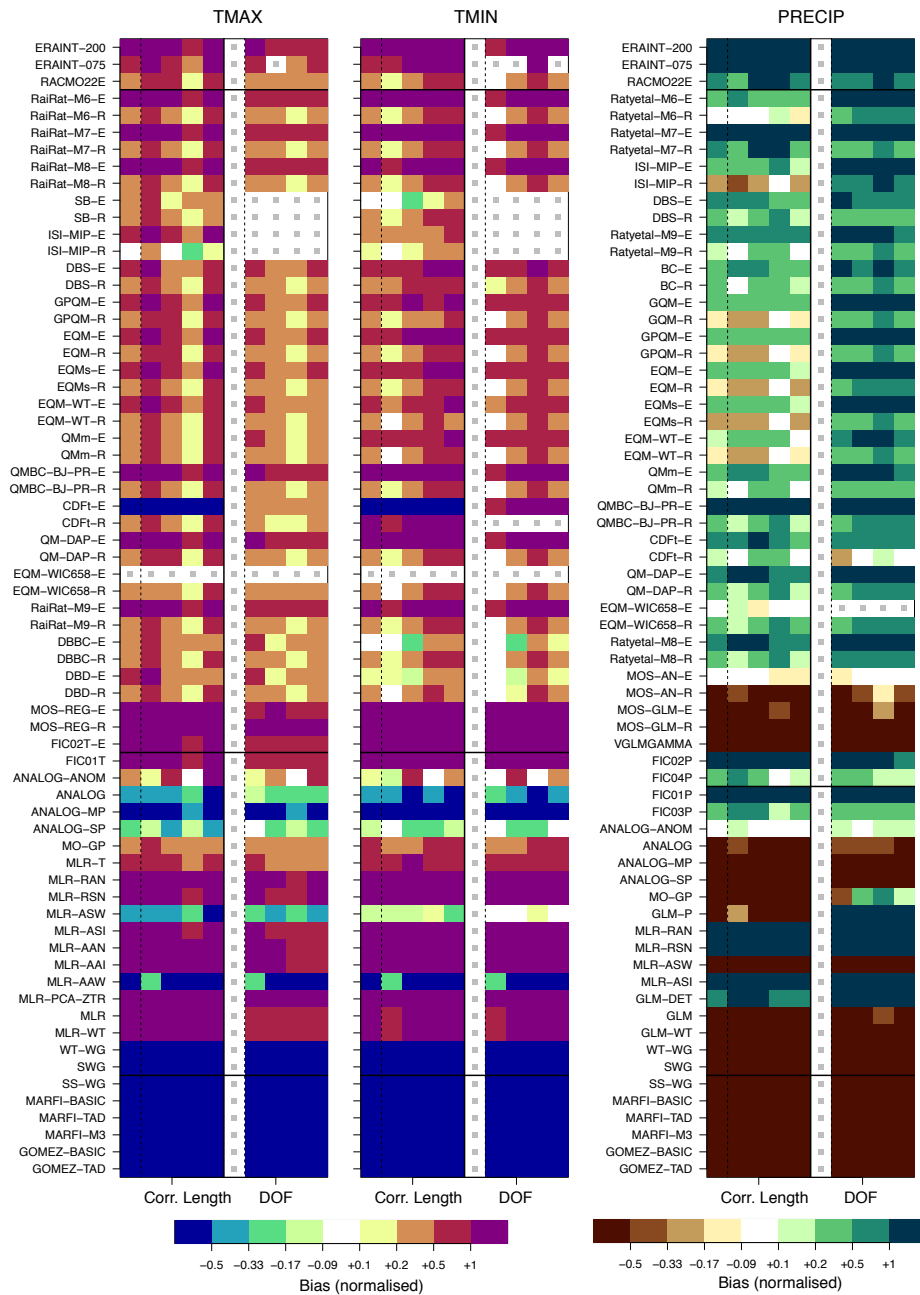


Figure 17: Relative biases in correlation length and independent spatial degrees of freedom (with sign swapped) based on the European network for daily maximum and minimum temperature, and precipitation. The columns indicate the seasons (annual, DJF, MAM, JJA, SON). For the degrees of freedom no annual values have been calculated.

830 bias in the degrees of freedom is not reduced as much. It was also shown that
831 MOS methods reduce the positive bias in the dependence for wet threshold
832 exceedance, but not for the exceedance of the 90th percentile of local daily pre-
833 cipitation. High-resolution, convection-permitting RCMs combined with MOS
834 might represent the spatial characteristics of heavy precipitation events consid-
835 erably better, but are still not widely used in climate change studies because
836 they are computationally expensive (Prein et al., 2015). The value added by the
837 regional model is still present after the MOS postprocessing (methods with suffix
838 '-R' perform better than those with suffix '-E'). For temperature the seasonal-
839 ity of the biases is similar for the raw model and for the MOS-corrected values.
840 The biases in correlation length and in the degrees of freedom are for minimum
841 temperature in general slightly higher than those for maximum temperature.

842 Fig. 17 and the specific findings in the main section also show that for all
843 predictand variables MOS methods perform in general better than PP methods,
844 however with some noteworthy exceptions. Deterministic PP methods that are
845 based on multiple linear regression and large-scale predictors tend to strongly
846 overestimate spatial correlations and also dependences of threshold exceedances,
847 while some other PP methods, for instance MO-GP and GLM-P, which use
848 local predictors, underestimate the joint variability between the stations, in
849 particular for precipitation. Given the different predictors used for different PP
850 methods it is possible that the results are strongly influenced by the predictor
851 choice rather than by the structure of the statistical model. Analog methods
852 yield, as expected, realistic spatial characteristics apart from sampling effects if
853 a common analog date is selected for all locations, whereas they underestimate
854 links between the stations if analogs are defined locally. In addition to the
855 analog methods the GLM-BN-DET method, which explicitly models spatial
856 dependence, performs very well with respect to the joint exceedance of the
857 local 90th percentile of daily precipitation, but somewhat underestimates the
858 joint exceedance of the wet-day threshold and of correlation lengths. Within the
859 set of PP methods analysed in our study multisite analog methods are thus the
860 only ones that are clearly suitable in applications where a realistic representation
861 of spatial variability is important. In climate change applications it needs to be
862 carefully checked however whether their use is justified, as potential changes of
863 the character of the analogs with respect to the predictor variable, and potential
864 new weather situation that are not well represented by the analogs may make
865 it difficult to capture the climate change signal. Furthermore, the temporal
866 sequence of the downscaled series might be unrealistic (Maraun et al., 2018).

867 The stochastic PP and MOS methods considered in the study yield time-
868 series that are too independent between the stations. There are two potential
869 contributions to this. First, the local variability that is explained by large-scale
870 predictors, and thus leads to links between locations, could be underestimated
871 due to the choice of statistical model and predictors. Second, the local noise is
872 independently added at different locations, and thus cannot include potential
873 links in the unexplained variability. The unconditional, local weather genera-
874 tors, which generate timeseries that are completely uncorrelated between the
875 locations, trivially fail to generate realistic spatial fields. Recently multisite
876 weather generators have been developed, and it has been demonstrated that
877 they can capture the spatial characteristics of precipitation at the catchment
878 scale well (e.g. Keller et al., 2015). If parameter changes in a future climate can
879 be credibly estimated, for instance by conditioning them on predictor variables,

880 such multisite weather generators can in principle be applied for climate change
881 studies.

882 As can be seen in Fig. 17 in most cases positive (negative) biases in the
883 correlation length are associated with negative (positive) biases in the degrees
884 of freedom, and the ranking of the magnitudes is similar. This might be ex-
885 pected as both measures are based on correlations and capture aspects of the
886 spatial complexity of the fields, with low (high) complexity likely to be associ-
887 ated with large (small) correlation lengths and a low (high) number of degrees
888 of freedom. However, there are some exceptions. For temperature the only
889 method for which the association is not found is CDFt-E, which as discussed
890 earlier might be due to technical problems with the method. The other excep-
891 tion are some of the MOS methods for precipitation, which have small negative
892 biases for the correlation lengths (see also Fig. 6) but also negative biases for
893 the degrees of freedom. This shows that although both measures usually yield
894 essentially the same information, subtleties in the correlation structure can exist
895 that lead to both biases having the same sign. This situation can occur because
896 the correlation lengths are dominated by station pairs with distances that lead
897 to correlations near the correlation threshold, whereas the degrees of freedom
898 are based on the entire correlation matrix. Although both approaches require
899 the calculation of the correlation matrix, calculating the degrees of freedom is
900 more straightforward because only the eigenvalue spectrum is required, whereas
901 determining the correlation lengths requires the calculation of correlations as
902 a function of distance, fitting of a smooth function, and involves a subjective
903 correlation threshold.

904 In summary we found that most PP downscaling methods yield unrealistic
905 spatial characteristics, regardless of whether large-scale or local predictors were
906 used, and therefore should not be applied for multisite downscaling if the spatial
907 characteristics of the results are relevant. The exception are multisite analog
908 methods and a method that explicitly models spatial dependence, which per-
909 formed well. The raw RCM clearly improves the skill compared to the driving
910 reanalysis. Adjusting the marginal distributions through MOS further reduces
911 biases in correlation lengths for precipitation and joint occurrence of wet days,
912 but does neither reduce the underestimation of complexity as measured by de-
913 grees of freedom, nor the substantial overestimation of the joint occurrence of
914 heavy precipitation events, while the improvements through the RCM are in
915 most cases retained. Whether the spatial characteristics of the output of these
916 methods is realistic enough for a given application needs to be carefully con-
917 sidered in each individual case. Moreover, a good performance in a perfect
918 predictor setup is no guarantee that the methods will perform well when driven
919 with GCM simulations for the present climate or that the climate change signal
920 is realistically represented (e.g. Maraun et al., 2017).

921 Despite the satisfying skill of some statistical downscaling methods, our re-
922 sults show that providing downscaled meteorological fields with realistic spatial
923 characteristics remains a challenge. In principle the common influence of predic-
924 tors in singlesite PP methods could lead to realistic spatial patterns, but in the
925 methods considered here it does not. The better skill of the RCM and of MOS
926 methods compared to most PP methods shows that explicit physical modelling
927 with local statistical post-processing is in general a better approach for obtaining
928 realistic spatial fields than deriving full spatial fields from large-scale predictors
929 (with the exceptions mentioned above). However none of the methods consid-

930 ered is able to produce output with a highly realistic spatial structure, including
931 the dependences for the exceedance of high precipitation thresholds. There is
932 thus still a clear need for increasing the resolution of RCMs used in climate
933 change studies, because the explicit physical modelling of small-scale processes
934 can be expected to improve the spatial characteristics of the raw model output
935 and of MOS-corrected fields, as well as lead to more realistic climate change
936 signals if regional processes affect climate change. Multisite weather generators
937 and multisite MOS have also the potential to yield realistic spatial fields, but
938 depend either on the assumption that the spatial dependence does not change
939 over time, or on ways to estimate and include changes in the dependence.

940 We note that the observation network used in VALUE is designed for val-
941 idation of a wide range of aspects of downscaling results, and not specifically
942 selected for the analysis of spatial variability. In particular the European net-
943 work, but also the German one, have station densities that do not well resolve
944 variability within small hydrological catchments. Thus similar studies with a
945 very high station density would be desirable. On very small scales subgrid vari-
946 ability becomes relevant for MOS methods and our results might not be directly
947 transferable because deterministic MOS approaches can be expected to lead to
948 too high dependences in cases where there is substantial subgrid variability
949 (Maraun, 2013).

950 As our intercomparison is based on an ensemble of opportunity of down-
951 scaling methods it would also be very useful to conduct future comparisons of
952 spatial aspects with a set of downscaling methods that does include all meth-
953 ods that are designed to represent spatial variability well. This should include
954 for instance the multisite weather generators and multisite MOS methods men-
955 tioned in the introduction. The evaluation of the former in different studies has
956 been inconclusive, while it has been positive for the latter, and a systematic
957 comparison using a common experimental setup would be very helpful for iden-
958 tifying suitable methods and for informing further method development. The
959 methods that explicitly model spatial dependence are more complex, more dif-
960 ficult to calibrate and apply, and more computationally expensive than most of
961 the methods used in our study, which is one of the main reasons they are not
962 frequently used and thus not included. The complexity of these methods also
963 means that they are not necessarily much easier to implement and apply than
964 high-resolution RCMs. Which combination of dynamical and statistical models
965 is best suited for a given application therefore needs careful consideration.

966
967

968 **Acknowledgements**

969 VALUE has been funded as EU COST Action ES1102. Participation of R. Huth
970 in VALUE was supported by the Ministry of Education, Youth, and Sports of
971 the Czech Republic under contract LD12059. JMG and SH acknowledge partial
972 funding from MULTI-SDM project (MINECO/FEDER, CGL2015-66583-R).

References

- 974 Arnaud, P., Bouvier, C., Cisneros, L. and Dominguez, R. (2002), ‘Influence of
975 rainfall spatial variability on flood prediction’, *J. Hydrol.* **260**(1), 216–230.
- 976 Ayar, P. V., Vrac, M., Bastin, S., Carreau, J., Déqué, M. and Gallardo, C.
977 (2016), ‘Intercomparison of statistical and dynamical downscaling models un-
978 der the EURO- and MED-CORDEX initiative framework: present climate
979 evaluations’, *Climate Dynamics* **46**(3-4), 1301–1329.
- 980 Bárdossy, A. and Pegram, G. (2012), ‘Multiscale spatial recorelation of RCM
981 precipitation to produce unbiased climate change scenarios over large areas
982 and small’, *Water Resour. Res.* **48**(9).
- 983 Boé, J., Terray, L., Habets, F. and Martin, E. (2006), ‘A simple statistical-
984 dynamical downscaling scheme based on weather types and conditional re-
985 sampling’, *J. Geophys. Res.-Atmos.* **111**(D23).
- 986 Bretherton, C. S., Widmann, M., Dymnikov, V. P., Wallace, J. M. and Blade, I.
987 (1999), ‘The effective number of spatial degrees of freedom of a time-varying
988 field’, *J. Clim.* **12**(7), 1990–2009.
- 989 Cannon, A. J. (2008), ‘Probabilistic multisite precipitation downscaling by an
990 expanded Bernoulli-Gamma density network’, *J. Hydrometeorol.* **9**(6), 1284–
991 1300.
- 992 Cannon, A. J. (2018), ‘Multivariate quantile mapping bias correction: an N-
993 dimensional probability density function transform for climate model simula-
994 tions of multiple variables’, *Clim. Dynam.* **50**(1-2), 31–49.
- 995 Cano, R., Sordo, C. and Gutiérrez, J. M. (2004), Applications of Bayesian
996 Networks in meteorology, in J. A. Gámez, S. Moral and A. Salmerón, eds,
997 ‘Advances in Bayesian Networks’, Springer Berlin Heidelberg, pp. 309–328.
- 998 Clark, M., Gangopadhyay, S., Hay, L., Rajagopalan, B. and Wilby, R. (2004),
999 ‘The Schaake shuffle: A method for reconstructing space-time variability in
1000 forecasted precipitation and temperature fields’, *J. Hydrometeorol.* **5**(1), 243–
1001 262.
- 1002 Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S.,
1003 Andrae, U., Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars,
1004 A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Delsol, C., Dragani, R.,
1005 Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Holm,
1006 E. V., Isaksen, L., Kallberg, P., Koehler, M., Matricardi, M., McNally, A. P.,
1007 Monge-Sanz, B. M., Morcrette, J. J., Park, B. K., Peubey, C., de Rosnay,
1008 P., Tavolato, C., Thepaut, J. N. and Vitart, F. (2011), ‘The ERA-Interim
1009 reanalysis: configuration and performance of the data assimilation system’,
1010 *Q. J. Roy. Meteor. Soc.* **137**(656, A), 553–597.
- 1011 Easterling, D. (1999), ‘Development of regional climate scenarios using a down-
1012 scaling approach’, *Climatic Change* **41**(3-4), 615–634.
- 1013 Ekstroem, M., Grose, M. R. and Whetton, P. H. (2015), ‘An appraisal of down-
1014 scaling methods used in climate change research’, *Wiley Interdisciplinary Re-
1015 views - Climate Change* **6**(3), 301–319.

- 1016 Ferraris, L., Gabellani, S., Reborá, N. and Provenzale, A. (2003), ‘A comparison
1017 of stochastic models for spatial rainfall downscaling’, *Water Resour. Res.*
1018 **39**(12).
- 1019 Frost, A., Charles, S. P., Timbal, B., Chiew, F., Mehrotra, R., Nguyen, K.,
1020 Chandler, R., McGregor, J., Fu, G., Kirono, D., Fernández, E. and Kent,
1021 M. (2011), ‘A comparison of multi-site daily rainfall downscaling techniques
1022 under Australian conditions’, *J. Hydrol.* **408**(1), 1–18.
- 1023 Gutiérrez, J., Maraun, D., Widmann, M., Huth, R., Hertig, E., Benestad, R.,
1024 Roessler, O., Wibig, J., Wilcke, R., Kotlarski, S., San-Martín, D., Herrera, S.,
1025 Bedia, J., Casanueva, A., Manzanás, R., Iturbide, M., Vrac, M., Dubrovsky,
1026 M., Ribalaygua, J., Pórtoles, J., Rätý, O., Räisänen, J., Hingray, B., Raynaud,
1027 D., Casado, M., Ramos, P., Zerenner, T., Turco, M., Bosshard, T., Štěpánek,
1028 P., Bartholy, J., Pongracz, R., Keller, D., Fischer, A., Cardoso, R., Soares, P.,
1029 Czernecki, B. and Pagé, C. (2018), ‘An intercomparison of a large ensemble
1030 of statistical downscaling methods over Europe: Results from the VALUE
1031 perfect predictor cross-validation experiment’, *Int. J. Climatol.* (in this issue).
- 1032 Gutmann, E., Pruitt, T., Clark, M. P., Brekke, L., Arnold, J. R., Raff, D. A.
1033 and Rasmussen, R. M. (2014), ‘An intercomparison of statistical downscaling
1034 methods used for water resource assessments in the United States’, *Water
1035 Resour. Res.* **50**(9), 7167–7186.
- 1036 Hannachi, A., Jolliffe, I. T. and Stephenson, D. B. (2007), ‘Empirical orthogonal
1037 functions and related techniques in atmospheric science: A review’, *Int. J.
1038 Climatol.* **27**(9), 1119–1152.
- 1039 Hengl, T. (2007), *A practical guide to geostatistical mapping of environmental
1040 variables*, European commission. Joint Research Centre. Publications Office,
1041 Luxembourg.
- 1042 Herdin, M., Czink, N., Ozcelik, H. and Bonek, E. (2005), Correlation matrix dis-
1043 tance, a meaningful measure for evaluation of non-stationary MIMO channels,
1044 in ‘Vehicular Technology Conference, 2005. VTC 2005-Spring. 2005 IEEE
1045 61st’, Vol. 1, IEEE, pp. 136–140.
- 1046 Hertig, E., Maraun, D., Bartholy, J., Pongracz, R., Vrac, M., Mares, I.,
1047 Gutiérrez, J. M., Wibig, J., Casanueva, A. and Soares, P. M. M. (2018), ‘Com-
1048 parison of statistical downscaling methods with respect to extreme events over
1049 Europe: Validation results from the perfect predictor experiment of the COST
1050 Action VALUE’, *Int. J. Climatol.* (in this issue).
- 1051 Hewitson, B. C., Daron, J., Crane, R. G., Zermoglio, M. F. and Jack, C.
1052 (2014), ‘Interrogating empirical-statistical downscaling’, *Climatic Change*
1053 **122**(4), 539–554.
- 1054 Hlinka, J., Hartman, D., Vejmelka, M., Novotna, D. and Palus, M. (2014), ‘Non-
1055 linear dependence and teleconnections in climate data: sources, relevance,
1056 nonstationarity’, *Clim. Dynam.* **42**(7-8), 1873–1886.
- 1057 Holzkämper, A., Calanca, P. and Fuhrer, J. (2012), ‘Statistical crop models:
1058 predicting the effects of temperature and precipitation changes’, *Clim. Res.*
1059 **51**, 11–21.

- 1060 Hu, Y., Maskey, S. and Uhlenbrook, S. (2013), ‘Downscaling daily precipitation
1061 over the Yellow River source region in China: a comparison of three statistical
1062 downscaling methods’, *Theor. Appl. Climatol.* **112**(3-4), 447–460.
- 1063 Hubert, L. and Arabie, P. (1985), ‘Comparing partitions’, *J. Classif.* **2**(2-
1064 3), 193–218.
- 1065 Huth, R. (2002), ‘Statistical downscaling of daily temperature in central Eu-
1066 rope’, *J. Clim.* **15**, 1731–1742.
- 1067 Huth, R., Kliegrova, S. and Metelka, L. (2008), ‘Non-linearity in statistical
1068 downscaling: does it bring an improvement for daily temperature in Europe?’,
1069 *Int. J. Climatol.* **28**(4), 465–477.
- 1070 Huth, R., Miksovsky, J., Stepanek, P., Belda, M., Farda, A., Chladova, Z.
1071 and Pisoft, P. (2015), ‘Comparative validation of statistical and dynamical
1072 downscaling models on a dense grid in central Europe: temperature’, *Theor.*
1073 *Appl. Climatol.* **120**(3-4), 533–553.
- 1074 Isotta, F. A., Vogel, R. and Frei, C. (2015), ‘Evaluation of European regional
1075 reanalyses and downscalings for precipitation in the Alpine region’, *Meteorol.*
1076 *Z.* **24**, 15–37.
- 1077 Keller, D., Fischer, A., Frei, C., Liniger, M., Appenzeller, C. and Knutti, R.
1078 (2015), ‘Implementation and validation of a Wilks-type multi-site daily pre-
1079 cipitation generator over a typical Alpine river catchment’, *Hydrol. Earth*
1080 *Syst. Sci.* **19**(5), 2163–2177.
- 1081 Kettle, H. and Thompson, R. (2004), ‘Statistical downscaling in European
1082 mountains: verification of reconstructed air temperature’, *Clim. Res.*
1083 **26**(2), 97–112.
- 1084 Kotlarski, S., Keuler, K., Christensen, O. B., Colette, A., Déqué, M., Gobiet,
1085 A., Goergen, K., Jacob, D., Lüthi, D., van Meijgaard, E., Nikulin, G., Schär,
1086 C., Teichmann, C., Vautard, R., Warrach-Sagi, K. and Wulfmeyer, V. (2014),
1087 ‘Regional climate modeling on European scales: a joint standard evaluation of
1088 the EURO-CORDEX RCM ensemble’, *Geosci. Model Dev.* **7**(4), 1297–1333.
- 1089 Machguth, H., Paul, F., Kotlarski, S. and Hoelzle, M. (2009), ‘Calculating dis-
1090 tributed glacier mass balance for the Swiss Alps from regional climate model
1091 output: a methodical description and interpretation of the results’, *J. Geo-*
1092 *phys. Res.* **114**.
- 1093 Mamalakis, A., Langousis, A., Deidda, R. and Marrocu, M. (2017), ‘A paramet-
1094 ric approach for simultaneous bias correction and high-resolution downscaling
1095 of climate model rainfall’, *Water Resour. Res.* **53**(3), 2149–2170.
- 1096 Maraun, D. (2013), ‘Bias correction, quantile mapping, and downscaling: Re-
1097 visiting the inflation issue’, *Journal of Climate* **26**(6), 2137–2143.
- 1098 Maraun, D., Huth, R., Gutiérrez, J. M., San-Martín, D., Dubrovsky, M., Fis-
1099 cher, A., Hertig, E., Soares, P. M. M., Bartholy, J., Pongrácz, R., Widmann,
1100 M., Casado, M. J., Ramos, P. and Bedia, J. (2018), ‘The VALUE perfect
1101 predictor experiment: evaluation of temporal variability’, *Int. J. Climatol.*
1102 (in this issue).

- 1103 Maraun, D., Shepherd, T. G., Widmann, M., Zappa, G., Walton, D., Gutiérrez,
1104 J. M., Hagemann, S., Richter, I., Soares, P. M., Hall, A. et al. (2017), ‘To-
1105 wards process-informed bias correction of climate change simulations’, *Nature*
1106 *Climate Change* **7**(11), 764.
- 1107 Maraun, D., Wetterhall, F., Ireson, A. M., Chandler, R. E., Kendon, E. J.,
1108 Widmann, M., Brienen, S., Rust, H. W., Sauter, T., Themessl, M., Venema,
1109 V. K. C., Chun, K. P., Goodess, C. M., Jones, R. G., Onof, C., Vrac, M.
1110 and Thiele-Eich, I. (2010), ‘Precipitation downscaling under climate change:
1111 Recent developments to bridge the gap between dynamical models and the
1112 end user’, *Rev. Geophys.* **48**(3). RG3003.
- 1113 Maraun, D. and Widmann, M. (2018), *Statistical Downscaling and Bias Cor-*
1114 *rection for Climate Research*, Cambridge University Press.
- 1115 Maraun, D., Widmann, M., Gutiérrez, J., Kotlarski, S., Chandler, R., Hertig,
1116 E., Wibig, J., Huth, R. and Wilcke, R. (2015), ‘VALUE: A framework to
1117 validate downscaling approaches for climate change studies’, *Earth’s Future*
1118 **3**(1), 1–14. 2014EF000259.
- 1119 Monestiez, P., Courault, D., Allard, D. and Ruget, F. (2001), ‘Spatial interpo-
1120 lation of air temperature using environmental context: application to a crop
1121 model’, *Environ. Ecol. Stat.* **8**, 297–309.
- 1122 Paschalis, A., Molnar, P., Fatichi, S. and Burlando, B. (2013), ‘A stochastic
1123 model for high-resolution space-time precipitation simulation’, *Water Resour.*
1124 *Res.* **49**(12), 8400–8417.
- 1125 Philipp, A., Della-Marta, P. M., Jacobeit, J., Fereday, D. R., Jones, P. D.,
1126 Moberg, A. and Wanner, H. (2007), ‘Long-term variability of daily North
1127 Atlantic-European pressure patterns since 1850 classified by simulated an-
1128 nealing clustering’, *J. Clim* **20**(16), 4065–4095.
- 1129 Pierce, D., Cayan, D., R. and Thrasher, B. (2014), ‘Statistical downscaling using
1130 localized constructed analogs (LOCA)’, *J. Hydrometeorol.* **15**(6), 2558–2585.
- 1131 Prein, A. F., Langhans, W., Fosser, G., Ferrone, A., Ban, N., Goergen, K.,
1132 Keller, M., Tölle, M., Gutjahr, O., Feser, F. et al. (2015), ‘A review on
1133 regional convection-permitting climate modeling: Demonstrations, prospects,
1134 and challenges’, *Reviews of geophysics* **53**(2), 323–361.
- 1135 R Core Team (2018), *R: A Language and Environment for Statistical Comput-*
1136 *ing*, R Foundation for Statistical Computing, Vienna, Austria.
1137 **URL:** <https://www.R-project.org/>
- 1138 Richman, M. (1986), ‘Rotation of principal components’, *J. Climatol.* **6**(3), 293–
1139 335.
- Santander Meteorology Group (2016), *R.VALUE: Climate data validation in*
the framework of the COST action VALUE. R package version 1.4-14.
URL: https://github.com/SantanderMetGroup/R_VALUE
- 1140 Santos, J. M. and Embrechts, M. (2009), On the use of the adjusted rand
1141 index as a metric for evaluating supervised classification, *in* ‘International
1142 Conference on Artificial Neural Networks’, Springer, pp. 175–184.

- 1143 Segond, M. L., Wheeler, H. S. and Onof, C. (2007), ‘The significance of spatial
1144 rainfall representation for flood runoff estimation: A numerical evaluation
1145 based on the Lee catchment, UK’, *J. Hydrol.* **347**(1), 116–131.
- 1146 Tank, A., Wijngaard, J., Können, G., Böhm, R., Demarée, G., Gocheva, A.,
1147 Mileta, M., Pashiardis, S., Hejkrlik, L., Kern-Hansen, C., Heino, R., Besse-
1148 moulin, P., Muller-Westermeier, G., Tzanakou, M., Szalai, S., Palsdottir, T.,
1149 Fitzgerald, D., Rubin, S., Capaldo, M., Maugeri, M., Leitass, A., Bukant-
1150 tis, A., Aberfeld, R., Van Engelen, A., Forland, E., Miletus, M., Coelho, F.,
1151 Mares, C., Razuvaev, V., Nieplova, E., Cegnar, T., Lopez, J., Dahlstrom,
1152 B., Moberg, A., Kirchhofer, W., Ceylan, A., Pachaliuk, O., Alexander, L.
1153 and Petrovic, P. (2002), ‘Daily dataset of 20th-century surface air tempera-
1154 ture and precipitation series for the European Climate Assessment’, *Int. J.*
1155 *Climatol.* **22**(12), 1441–1453.
- 1156 Trambauer, P., Werner, M., Winsemius, H. C., Maskey, S., Dutra, E. and
1157 Uhlenbrook, S. (2015), ‘Hydrological drought forecasting and skill assess-
1158 ment for the Limpopo river basin, Southern Africa’, *Hydrol. Earth Syst. Sci.*
1159 **19**(4), 1695–1711.
- 1160 van Meijgaard, E., van Ulf, L. H., van de Berg, W. J., Bosveld, F. C., van den
1161 Hurk, B. J. J. M., Lenderink, G. and Siebesma, A. P. (2008), The KNMI
1162 regional atmospheric climate model RACMO version 2.1, Technical Report
1163 302, Royal Dutch Meteorological Institute, KNMI, Postbus 201, 3730 AE, De
1164 Bilt, The Netherlands.
- 1165 Viviroli, D., Zappa, M., Schwanbeck, J., Gurtz, J. and Weingartner, R. (2009),
1166 ‘Continuous simulation for flood estimation in ungauged mesoscale catch-
1167 ments of Switzerland - Part I: modelling framework and calibration results’,
1168 *J. Hydrol* **377**, 191–207.
- 1169 Voisin, N., Pappenberger, F., Lettenmaier, D. P., Buizza, R. and Schaake, J. C.
1170 (2011), ‘Application of a medium-range global hydrologic probabilistic fore-
1171 cast scheme to the Ohio River Basin’, *Weather Forecast.* **26**(4), 425–446.
- 1172 Vrac, M. (2018), ‘Multivariate bias adjustment of high-dimensional climate sim-
1173 ulations: The “Rank Resampling for Distributions and Dependences” (R2D2)
1174 bias correction’, *Hydrol. Earth Syst. Sci. Discuss.* **2018**, 1–33.
- 1175 Vrac, M. and Friederichs, P. (2015), ‘Multivariate-intervariable, spatial and
1176 temporal-bias correction’, *J. Clim.* **28**, 218–237.
- 1177 Wilks, D. S. (2012), ‘Stochastic weather generators for climate-change down-
1178 scaling, part II: multivariable and spatially coherent multisite downscaling:
1179 Stochastic weather generators for climate-change downscaling’, *Wiley Inter-*
1180 *disciplinary Reviews: Climate Change* **3**(3), 267–278.